

Broadview  
www.broadview.com.cn

Regular Expression  
yet another introduction

# 正则指引

余晟 著



电子工业出版社  
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

# 正 则 指 引

余晟 著

電子工業出版社

Publishing House of Electronics Industry

北京 · BEIJING

## 内 容 简 介

本书针对作者在开发中遇到的实际问题,以及其他开发人员咨询的问题,总结出一套使用正则表达式解题的办法,并通过具体的例子指导读者拆解、分析问题。全书分为三大部分:第一部分主要讲解正则表达式的基础知识,涵盖了常见正则表达式中的各种功能和结构;第二部分主要讲解关于正则表达式的更深入的知识,详细探讨了编码问题、匹配原理、解题思路;第三部分将之前介绍的各种知识落实到6种常用语言.NET、Java、JavaScript、PHP、Python、Ruby中,不但详细介绍了语言中正则表达式的用法,更点明了版本之间的细微差异,既可以作为专门学习的教材,也可以作为有用的参考手册。

本书适合经常需要进行文本处理(比如日志分析或网络运维)的技术人员、熟悉常用开发语言的程序员,以及已经对正则表达式有一定了解的读者阅读。

未经许可,不得以任何方式复制或抄袭本书之部分或全部内容。  
版权所有,侵权必究。

### 图书在版编目(CIP)数据

正则指引 / 余晟著. — 北京: 电子工业出版社, 2012.5  
ISBN 978-7-121-16551-1

I. ①正… II. ①余… III. ①正则表达式 IV. ①TP301.2

中国版本图书馆CIP数据核字(2012)第046409号

策划编辑: 张月萍

责任编辑: 葛 娜

印 刷:

装 订: 中国电影出版社印刷厂

出版发行: 电子工业出版社

北京市海淀区万寿路173信箱 邮编: 100036

开 本: 787×980 1/16 印张: 21 字数: 450千字

印 次: 2012年5月第1次印刷

印 次: 4000册 定价: 58.00元

凡所购买电子工业出版社图书有缺损问题,请向购买书店调换。若书店售缺,请与本社发行部联系,联系及邮购电话:(010) 88254888。

质量投诉请发邮件至 [zlts@phei.com.cn](mailto:zlts@phei.com.cn), 盗版侵权举报请发邮件到 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

服务热线:(010) 88258888。

# 引子：关于正则表达式……

**正则表达式**这个名字看起来总有点古怪，概念似乎也不简单，甚至需要用一整本书来讲解；可是，它到底是什么呢？

身为技术人员，我相信你总会与字符串打交道，相应地，各种语言也都提供了与字符串有关的函数。不妨先看看下面几个问题，字符串函数是如何解决的（下面的代码使用 **Python** 语言，它很直观，正文里有基础的介绍。现在，你只需要知道 **def** 是定义函数的关键词即可）。

## 1. 判断字符 **ch** 是否数字字符

```
def isDigit(ch) :  
    return ch == "0" or ch == "1" ..... or ch == "9"
```

2. 判断字符串 **str** 是否电话号码( 为简单起见, 现在只考虑固定电话号码, 也就是长度在 7~8 位之间的数字字符串, 且第一位不为 0 )

```
def isPhoneNum (str) :  
    if len(str) >= 7 and len(str) <= 8 and str[0] != "0" :  
        for ch in str :  
            if not isDigit(ch) :  
                return false  
        return true  
    return false
```

任务本身并没有增加太多，但是程序复杂了很多倍；如果你不这样看，那么，来个更复杂的。

## 3. 找出一段文本中所有的电话号码

最直观的办法是，在字符串中的每个位置截取 7~8 个字符，调用之前的 **isPhoneNum()**。这么做看起来没问题，只是效率太低。

当然，很容易就可以做点改进，只在“当前字符为数字字符”的情况下调用 **isPhoneNum()**。这样效率倒是改进了，但是还有问题没有解决：要求找到的是长度大于等于 7 个字符，小于等于 8 个字符的“数字字符串”，而不是“子字符串”——也就是说，假如数字字符串是 **64240000**，需要将它找出来；如果数字字符串是 **13800138000**，则需要忽略它，以及其中的任何子串（比

如 13800138, 00138000)。

所以，用 `isPhoneNum()` 找出字符串之后，还需要保证它之前的字符不是数字字符，之后的字符也不是数字字符。看起来很简单，但为了避免越界错误，又需要判断：如果当前字符是整段文本的第一个字符，则不需要判断之前的字符，因为它不存在；同样，如果找出的字符串在整段文本的末尾，则不需要判断之后的字符，因为它同样不存在……

到现在为止，即便只是找到最简单的固定电话号码，程序也非常复杂，难以维护。如果要查找的是形式更多变的文本，比如带区号的电话号码（021-64240000 或者 03718888888）、手机号码（13800138000 或者 +8613800138000 或者 013800138000），程序更是不可想象，更不用说文件路径名、URL 地址、电子邮件地址了！然而，日常开发中我们又确实经常需要面对这类任务，有什么更好的办法呢？

正则表达式就是解决这类问题的万灵药。虽然许多人有点看不起它，觉得不入流，科班教材里也不会花太多篇幅来介绍它，但它确实是解决问题的利器——之前提到的三个例子，用正则表达式都可以轻松解决。

### 1. 判断字符 ch 是否数字字符

```
def isDigit(ch) :  
    return re.search(ch, "[0-9]") != None
```

看起来很复杂，其实并不复杂：这里真正要关心的就是正则表达式 `[0-9]`，它表示“从 0 到 9 之间的任意字符”，很形象吧？`re.search()` 是正则表达式运算函数，它判断 `ch` 能否由正则表达式 `[0-9]` 匹配，可以则返回一个结果，否则返回 `None`（这些细节正文中会讲到）。

### 2. 判断字符串 str 是否电话号码

```
def isPhoneNum(str) :  
    return re.search(str, "[1-9][0-9]{6,7}") != None
```

这个正则表达式最开始是 `[1-9]`，表示第一个字符必须是 1~9 之间的数字字符；之后是 `[0-9]{6,7}`，表示长度在 6 和 7 之间，由 0~9 之间的数字字符组成的字符串（两部分加起来，整个字符串的长度在 7 和 8 之间）。要解决的问题复杂了，正则表达式仍然直观形象。

### 3. 找出一段文本中所有的固定电话号码

```
def findNumStr(str) :  
    return re.findall(str, '(?![0-9])[1-9][0-9]{6,7}(?![0-9])')
```

这个正则表达式之前多出了 `(?![0-9])`，表示“之前不能是 [0-9]”；之后多出了 `(?![0-9])`，表示“之后不能是 [0-9]”。虽然稍微复杂点，但意思明确，而且不难理解。`re.findall()` 的意思也很明显：找到所有这样的字符串。

可以想象，循着这种思路，更复杂的电话号码、手机号码等任务都不难解决。更重要的是，

之前需要许多行语句才能完成的任务，现在基本上只需要一个正则表达式，一条语句就可以完成。正因为如此，不少人虽然认为正则表达式不够花哨、漂亮，却不得不承认它是一种“匕首应用”——匕首，没有十八般兵刃那么大方，关键时候却不可或缺，所以值得花时间练练。同样，正则表达式虽然不能用来显摆，但总有派得上用场的地方，花时间练练绝不是坏事。即便你的工作不是纯粹的文本处理（比如日志分析），也总会有用到正则表达式的地方（比如查找和修改源代码），所以我希望，这本书能陪伴你练出一身正则表达式的好功夫，在关键场合能亮出趁手的工具。

最后，为了尊重传统教科书的习惯，附上正则表达式的“科班史”：

正则表达式发源于与计算机密切相关的两个领域：计算理论和形式语言。20 世纪 40 年代，两位神经生理学家 **Warren McCulloch** 和 **Walter Pitts** 研究出一种数学方式来描述神经网络的办法，它们把神经系统中的神经元描述成小而简单的自动控制单元。1956 年，数学家 **Stephen Cole Kleene** 在他们研究的基础上，发表了一篇名为“神经网络事件的表示法”的论文，在其中，他采用了一些称之为“正则集合（regular set）”的数学符号来描述神经网络模型。

之后，UNIX 的主要发明人 **Ken Thompson** 将这个符号系统引入了文本编辑器 **QED**（意思是“在文本中搜索某种模式”），正则表达式由此也进入了计算机世界。随后 **Ken Thompson** 又将正则表达式引入了 UNIX 下的文本编辑器 **ed**，**ed** 最终演化为大家熟悉的 **grep**（**grep** 得名自 **ed** 编辑器中的正则表达式搜索命令 **g/re/p**，其中的 **re** 表示“正则表达式”）。

# 返璞归真——评《正则指引》

第一次接触正则表达式，是 2000 年我在西安一家公司使用 Perl 做网站开发时。之前我在工作中只使用过标准的 C 语言，Perl 这门编程语言的强大表达能力，令我印象极为深刻。Perl 的力量，除了语言本身的设计之外，很大程度上来自于它对正则表达式的完美支持。当时我们开发了一个网上商城的应用，允许很多商家在这里开店，可以选择一些不同的样式模板。我很快发现，使用 Perl+正则表达式是开发这类应用的利器。我们只花了大约一个月的时间，就完成了网站核心功能的开发。那时候我意识到，使用正则表达式是聪明人写程序的方法（没说我是聪明人，但是我非常希望与那些聪明人为伍），可以极大地提高代码的重用度和执行效率。如果完全不使用正则表达式，代码量会增加数倍甚至十倍。

后来因为一些原因，我告别了 Perl。在之后的工作中，我使用过 Java、JavaScript、Ruby 等编程语言。我发现这些语言对于正则表达式的支持，没有一个能够超越 Perl。Java 这种所谓的“工业主流编程语言”，一直到 2002 年 JDK 1.4 推出时，才正式把对正则表达式的支持加入到核心类库。因为长期缺乏对正则表达式的原生支持，以及语言本身表达能力欠缺，使用 Java 来做大量的文本处理，感觉非常笨拙，完全没有使用 Perl 那种指哪打哪的快感。直到 2007 年我发现了另一个更好的 Perl 语言——Ruby，才重新找回了 2000 年 Perl 带给我的编程快感。

因为我的工作主要是做 Web 开发，大量的时间花在与 HTML/CSS/JavaScript 以及关系数据库打交道上。在这里并没有很高深的算法，只有大量繁重的文本处理。难以想象，如果没有正则表达式，我们的开发将会是何等原始。

除了 Web 开发领域，需要实现大量自动化功能的一些领域，例如运维领域和自动化测试领域，也是正则表达式大显身手的地方。无论使用稍显简陋的 sed/awk 还是更高级的 Perl/Python/Ruby，实现自动化功能，都必须依赖大量的正则表达式。

自从面向对象时髦起来之后，甚至一度出现了面向对象万能论，有人试图用 MDA 和可执行的 UML 来解决一切编程问题。但是我一直认为面向对象只解决了软件开发的一小部分问题，而且是宏观方面的问题。正则表达式解决的问题，是面向对象无能为力的一些微观方面的问题。在这里不需要坐而论道的方法论争论，需要的是刺刀见红的肉搏战。这些问题即使使用完全面向对象的方式能够解决，也会是很笨拙的。如果用物理学来比喻，面向对象是“广义相对论”，而正

则表达式则是“量子力学”。

正则表达式已经成为了现代编程语言的基础模块，现在很难找到一种不支持正则表达式的编程语言。除了编程语言外，在很多工具软件，例如文本编辑器（Vi、Emacs、UltraEdit）、Web 服务器（Apache、Nginx）之中都能找到正则表达式的身影。

余晟老师是我的朋友，我对他印象最为深刻的是他对于技术工作的严谨态度。“格物致知”是中国传统儒家学派所追求的一种道德修养，也是一种境界。余老师是我的朋友中最接近“格物致知”这种境界的一位。我虽然从未精通过任何一门技术，但是很喜欢结交余老师这样的朋友。

余老师潜心编著的这本《正则指引》深入浅出，将正则表达式的由来和分支娓娓道来。阅读这本书，我仿佛回到了 11 年前做 Perl 程序员时的快乐时光。国内很多程序员的一个通病是好高骛远，像《正则指引》这样一本详细讲解基础知识的书未必会有很好的销路。但是等你做过很多年开发之后，你会发现，对你最有价值的，正是这些基础知识和工具。软件开发的“道”，正是隐藏在这些看起来不起眼的基础知识和工具之中的。

李锟

2011 年 11 月 25 日



# 前言

提到正则表达式，许多人很有点不屑一顾：这东西，不登大雅之堂，再说也不是总要用到，何必专门花时间学习？

没错，正则表达式并不“总要用到”，但到了需要的场合用不上，往往产生“一分钱难倒英雄汉”的尴尬。经常需要处理文本的程序员自然会知道正则表达式的价值，其他的程序员如果不会正则表达式，即便开发的领域与文本处理没什么关系，也难免“躺着中枪”的命运——前几天我遇到一个问题，将一行长长的地址拆分成多行，负责这部分的程序员的日常工作只是制作 PDF 而已，拆分地址是很“边缘”的功能，但不会正则表达式就无法准确折行（一般需要在标点符号出现的地方折行，而不能只在空白字符处折行，但是不同语言中的标点符号各有不同），结果一筹莫展；相反，如果了解正则表达式，就可以很容易地处理各种语言中的标点字符。

以我的开发经验来看，专门花点时间掌握正则表达式，确实是非常有必要的。目前可以见到的关于正则表达式的书籍和资料有不少，但又各有不足。

在互联网上，流传着一些编程语言的正则文档和《30 分钟教会你正则表达式》之类的帖子。这类资料的好处是简单直接，查到了，如果有现成的例子，而且适用于自己的语言，则可以直接拿来用；然而，其坏处也是简单直接，因为缺乏背后原理的讲解，如果找不到现成的例子，或者找不到能在自己所使用语言中行得通的例子（需知道，同样的正则表达式并不能直接套用到不同的语言中），则束手无策。

在正式的出版领域，已经有《精通正则表达式》、《正则表达式必知必会》之类的书籍出版，尤其是前者，堪称关于正则表达式的经典著作，如果想认真学习正则表达式，这类书籍是必须阅读的。但这类书籍也有一个弱点，即都是由英文版本翻译而来的，更多地侧重英文文本的处理，身为中文世界的开发人员，我们经常需要处理中文文本——对于处理英文之外的字符，正则表达式已经提供了足够丰富的功能，但如何用对、用好这些功能，资料却很匮乏。

我经常需要给人讲解正则表达式的相关知识，时常惋惜的是，开发人员为这些问题所困扰；正因为如此，本书的写作动机便是着力弥补现有资料的缺陷。

相对于正则文档和速成教学帖子，本书深入讲解了匹配背后的原理，往往会举一反三，告诉

读者，这里为何这样写，如果改成其他形式，会造成什么结构；并且，集中讲解和比较了多种语言中正则表达式用法的异同，方便读者把现成的正则表达式“移植”到自己的工作环境中。

相对于《精通正则表达式》等正式的书籍，本书辟出专门的内容讲解语言和编码，告诉读者如何设定编码，如何正确处理中文等字符。另外，本书还涵盖了.NET、Java、JavaScript、PHP、Python、Ruby 六种常用语言，对每种语言给出专门章节，不但详细介绍了语言中正则表达式的用法，更点明了版本之间的细微差异，不但可以作为专门学习的教材，还可以成为有用的参考手册。

## 本书结构

本书可以分为三大部分。

第一部分主要讲解正则表达式的基础知识，覆盖常见正则表达式中的各种功能和结构。看完前面 3 章，就可以基本弄明白现在流行的各种正则表达式；尤其是如果你之前有一些经验，会觉得阅读起来并不困难。但是我也希望读者不要忽略其他的内容，断言和匹配模式现在已经是正则表达式的“标准配备”了，而且确实可以派上大用场，所以第 4 章和第 5 章的内容，即便不是很熟悉，阅读起来可能有一些麻烦，也不应该忽略。最后的第 6 章，则厘清了正则表达式在使用中的若干疑惑，了解它们，你就可以相对自由地在正则表达式的世界里行走了。

第二部分主要讲解关于正则表达式的更深入的知识，这一部分用 3 章的内容，详细探讨了编码问题、匹配原理、解题思路。这部分内容更抽象，需要多花一点时间来阅读和理解，但是它们确实可以帮你在正则表达式的世界里登堂入室，脱离“术”的层面，掌握万变不离其宗的“道”。

第三部分的作用是接地气，将之前介绍的各种知识落实到六种常用语言.NET、Java、JavaScript、PHP、Python、Ruby 中来。每一章的开头有正则功能列表，其中的功能都对应着前面部分的讲解，这些功能的具体应用实例，以及不同版本之间的差异，则在章节中详细讲解，每一章的最后还给出了常见任务的示例代码，方便日后查询。在最后，第 16 章简要介绍了正则表达式在 Linux 下常用工具 vi、grep、awk、sed 中的使用，并通过一个实际的例子将这几种工具串起来，对比说明了它们适合解决的问题。

在本书的最后提供了用作参考的 3 个附录。

附录 A 是正则表达式的常用功能在不同语言中的比对，希望能给需要在多种语言中使用正则表达式或者移植正则表达式的读者提供一份有用的参考；附录 B 给出了若干常见的正则表达式，比如匹配邮政编码、身份证号、手机号、QQ 号、电子邮件地址等，希望能成为常见问题的“速查手册”；附录 C 列出了常用正则表达式的工具和资源，方便大家调试自己的正则表达式，以及继续深入学习。

## 本书读者

本书适合以下几类读者。

经常需要进行文本处理（比如日志分析或网络运维）的技术人员。这些读者或许已经熟悉了正则表达式的基本用法，但面对日益复杂化和海量化的数据，阅读本书可以帮助你更准确、更高效地处理文本，提升自己工作的价值。

熟悉常用开发语言的程序员。虽然这些读者不需要专职进行文本处理，但源代码和许多数据其实也是文本，如果不会正则表达式，在偶然遇到处理源代码或文本数据的任务时，往往会产生躺着中枪的无力感。本书第三部分可以帮你迅速找到有关的例子，并落实在自己的编程语言中。当然前两部分也非常有必要，因为它们可以帮你夯实基础。

对正则表达式已经有一定了解的读者。这些读者虽然能用正则表达式解决常见的任务，但未必了解正则表达式的编码问题、匹配原理、解题思路，仔细阅读本书的第二部分，可以深化完善对正则表达式的理解；而第三部分详细比较了使用正则表达式时各种语言，以及同一种语言中各种版本的差异。所有这一切，应该可以让你对正则表达式的掌握更上一层楼。

## 致谢

一本书的完成，必然离不开众多人的帮忙。

首先需要感谢的是周筠老师和徐定翔、卢鹤翔两位编辑，他们在我写作的最初阶段做了大量细心、耐心的工作，完全可以说，没有他们的这些工作，我就不会有写作这本书的念头，或者坚持写完的动力。

然后要感谢的是电子工业出版社的杨福平副总编和张月萍编辑，没有他们的关照和辛劳工作，这本书的出版定然会遇到更多的困难。

感谢我的朋友霍炬和韩磊，虽然我之前阅读过《精通正则表达式》，但与翻译和写作结缘，他们给了我莫大的帮助，有了这个契机，才有了现在的《正则指引》。尤其值得一提的是霍炬的夫人西乔，精心手绘了这本书的封面，在这里表示诚挚的谢意。

感谢我曾工作过的盛大创新院以及创新院的各位同事（李骏、郝培强、庄表伟、丁宇、许式伟、莫华枫、李道兵、赵劼、樊一鹏、张一宁等），创新院给大家宽松自由的工作环境，与各位同事的讨论加深了我对正则表达式的理解，也为我提供了许多例子。

感谢张东亮、陆亦斌、孙勇、叶劲峰等各位朋友，愿意拨冗阅读本书的草稿，并提出了大量专业的意见。

感谢何源、陈钢、贺钧、陈驰等读者，试读本书之后提出了大量的宝贵意见，在最后关头打

消了我心中的许多忐忑。

在更早之前，我的父母从小就鼓励我研究和了解各种科学原理（“玩也要动脑筋”），我之所以有兴趣探究正则表达式背后的世界，而不满足于“够用/凑合”，归源都是受益于这种思维行为习惯。此外，在中小学阶段，我的语文老师罗碧玉、郭志鸿、易玺铭培养了我对于文字的兴趣，在大学阶段，东北师范大学文学学院的王确老师给了我这个理科生非常多的帮助和指引。对各位师长，在此一并表示感谢，能遇到你们是我的幸运。

最后还需要感谢许多为这本书做出过贡献的人，你们的名字我可能暂时无法记起，或者无法一一罗列，但我会心中存留对你们的谢意。

# 目 录

## 第 一 部 分

第 1 章 字符组 .....	2
1.1 普通字符组.....	2
1.2 关于 Python 的基础知识 .....	4
1.3 普通字符组（续） .....	6
1.4 元字符与转义.....	8
1.5 排除型字符组.....	10
1.6 字符组简记法.....	12
1.7 字符组运算.....	14
1.8 POSIX 字符组 .....	15
第 2 章 量词 .....	17
2.1 一般形式 .....	17
2.2 常用量词 .....	18
2.3 数据提取 .....	21
2.4 点号 .....	23
2.5 滥用点号的问题.....	23
2.6 忽略优先量词.....	26
2.7 转义 .....	31
第 3 章 括号 .....	33
3.1 分组 .....	33
3.2 多选结构 .....	39
3.3 引用分组 .....	44
3.3.1 反向引用.....	48
3.3.2 各种引用的记法.....	50
3.3.3 命名分组.....	53
3.4 非捕获分组.....	54

3.5	补充 .....	55
3.5.1	转义.....	55
3.5.2	URL Rewrite .....	56
3.5.3	一个例子.....	58
第 4 章	断言 .....	59
4.1	单词边界 .....	59
4.2	行起始/结束位置.....	61
4.3	环视 .....	68
4.4	补充 .....	74
4.4.1	环视的价值.....	74
4.4.2	环视与分组编号 .....	74
4.4.3	环视的支持程度.....	75
4.4.4	环视的组合.....	77
4.4.5	断言和反向引用之间的关系 .....	79
第 5 章	匹配模式 .....	81
5.1	不区分大小写模式.....	81
5.1.1	模式的指定方式 .....	82
5.2	单行模式 .....	84
5.3	多行模式 .....	85
5.4	注释模式 .....	87
5.5	补充 .....	88
5.5.1	更多的模式.....	88
5.5.2	修饰符的作用范围.....	89
5.5.3	失效修饰符.....	90
5.5.4	模式与反向引用 .....	90
5.5.5	冲突策略.....	91
5.5.6	哪种方式更好 .....	92
第 6 章	其他 .....	93
6.1	转义 .....	93
6.1.1	字符串转义与正则转义.....	93
6.1.2	元字符的转义.....	97
6.1.3	彻底消除元字符的特殊含义 .....	99
6.1.4	字符组中的转义.....	101
6.2	正则表达式的处理形式.....	101
6.2.1	函数式处理.....	102

6.2.2	面向对象式处理 .....	102
6.2.3	比较 .....	103
6.2.4	线程安全性 .....	104
6.3	表达式中的优先级 .....	106

## 第 二 部 分

第 7 章	Unicode .....	110
7.1	关于编码 .....	110
7.2	推荐使用 Unicode 编码 .....	111
7.3	Unicode 匹配规则 .....	115
7.4	单词边界 .....	117
7.5	码值 .....	119
7.6	Unicode 属性 .....	121
7.6.1	Unicode Property .....	121
7.6.2	Unicode Block .....	122
7.6.3	Unicode Script .....	123
7.7	Unicode 属性列表 .....	123
7.7.1	Unicode Property .....	123
7.7.2	Unicode Block .....	125
7.7.3	Unicode Script .....	128
7.8	POSIX 字符组 .....	129
第 8 章	匹配原理 .....	130
8.1	有穷自动机 .....	130
8.2	正则表达式的匹配过程 .....	131
8.3	回溯 .....	134
8.4	NFA 和 DFA .....	136
第 9 章	常见问题的解决思路 .....	138
9.1	关于元素的三种逻辑 .....	138
9.1.1	必须出现 .....	139
9.1.2	可能出现 .....	139
9.1.3	不能出现 .....	140
9.2	正则表达式的常见操作 .....	142
9.2.1	提取 .....	142
9.2.2	验证 .....	148

9.2.3	替换.....	152
9.2.4	切分.....	157
9.3	正则表达式的优化建议.....	159
9.3.1	使用缓存.....	159
9.3.2	尽量准确地表达意图.....	160
9.3.3	避免重复匹配.....	160
9.3.4	独立出文本和锚点.....	161
9.4	别过分依赖正则表达式.....	162
9.4.1	彻底放弃字符串操作.....	162
9.4.2	思维定势.....	163
9.4.3	正则表达式可以匹配各种文本.....	164

## 第 三 部 分

第 10 章	.NET.....	168
10.1	预备知识.....	168
10.2	正则功能详解.....	169
10.2.1	列表.....	169
10.2.2	字符组.....	170
10.2.3	Unicode 属性.....	170
10.2.4	字符组简记法.....	171
10.2.5	单词边界.....	171
10.2.6	行起始/结束位置.....	172
10.2.7	环视.....	173
10.2.8	匹配模式.....	173
10.2.9	捕获分组的引用.....	174
10.3	正则 API 简介.....	175
10.3.1	Regex.....	175
10.3.2	Match.....	179
10.4	常用操作示例.....	180
10.4.1	验证.....	180
10.4.2	提取.....	180
10.4.3	替换.....	181
10.4.4	切分.....	182
第 11 章	Java.....	183
11.1	预备知识.....	183



11.2	正则功能详解.....	184
11.2.1	列表.....	184
11.2.2	字符组.....	184
11.2.3	Unicode 属性 .....	186
11.2.4	字符组简记法.....	186
11.2.5	单词边界.....	186
11.2.6	行起始/结束位置.....	187
11.2.7	环视.....	188
11.2.8	匹配模式.....	188
11.2.9	纯文本模式.....	189
11.2.10	捕获分组的引用 .....	189
11.3	正则 API 简介 .....	189
11.3.1	Pattern .....	190
11.3.2	Matcher .....	192
11.3.3	String.....	194
11.4	常用操作示例.....	195
11.4.1	验证.....	195
11.4.2	提取.....	196
11.4.3	替换.....	196
11.4.4	切分.....	197
第 12 章	JavaScript .....	198
12.1	预备知识.....	198
12.2	正则功能详解.....	199
12.2.1	列表.....	199
12.2.2	字符组.....	199
12.2.3	字符组简记法.....	200
12.2.4	单词边界.....	200
12.2.5	行起始/结束位置.....	201
12.2.6	环视.....	201
12.2.7	匹配模式.....	202
12.2.8	捕获分组的引用.....	203
12.3	正则 API 简介 .....	203
12.3.1	RegExp.....	203
12.3.2	String.....	207
12.4	常用操作示例.....	210

12.4.1	验证.....	210
12.4.2	提取.....	210
12.4.3	替换.....	211
12.4.4	切分.....	211
12.5	关于 <code>ActionScript</code> .....	211
12.5.1	<code>RegExp</code> .....	211
12.5.2	匹配规则.....	212
12.5.3	匹配模式.....	212
12.5.4	正则 <code>API</code> .....	212
第 13 章	<code>PHP</code> .....	213
13.1	预备知识.....	213
13.2	正则功能详解.....	215
13.2.1	列表.....	215
13.2.2	字符组.....	216
13.2.3	<code>Unicode</code> 属性 .....	217
13.2.4	字符组简记法.....	217
13.2.5	单词边界.....	217
13.2.6	行起始/结束位置.....	218
13.2.7	环视.....	219
13.2.8	匹配模式.....	219
13.2.9	纯文本模式.....	220
13.2.10	捕获分组的引用.....	220
13.3	正则 <code>API</code> 简介 .....	221
13.3.1	<code>PREG</code> 常量说明.....	221
13.3.2	<code>preg_quote</code> .....	222
13.3.3	<code>preg_grep</code> .....	223
13.3.4	<code>preg_match</code> .....	223
13.3.5	<code>preg_match_all</code> .....	225
13.3.6	<code>preg_last_error</code> .....	227
13.3.7	<code>preg_replace</code> .....	227
13.3.8	<code>preg_replace_callback</code> .....	227
13.3.9	<code>preg_filter</code> .....	228
13.3.10	<code>preg_split</code> .....	229
13.4	常见的正则操作举例.....	230
13.4.1	验证.....	230

13.4.2	提取.....	230
13.4.3	替换.....	231
13.4.4	切分.....	232
第 14 章	Python.....	233
14.1	预备知识.....	233
14.2	正则功能详解.....	234
14.2.1	列表.....	234
14.2.2	字符组.....	235
14.2.3	Unicode 属性 .....	236
14.2.4	字符组简记法.....	236
14.2.5	单词边界.....	238
14.2.6	行起始/结束位置.....	239
14.2.7	环视.....	239
14.2.8	匹配模式.....	240
14.2.9	捕获分组的引用.....	240
14.3	正则 API 简介 .....	241
14.3.1	RegexObject.....	241
14.3.2	re.compile(regex[, flags]) .....	243
14.3.3	re.search(pattern, string[, flags]) .....	243
14.3.4	MatchObject.....	243
14.3.5	re.match(pattern, string[, flags]).....	244
14.3.6	re.findall(pattern, sting[, flags]) .....	245
14.3.7	re.finditer(pattern, string[, flags]).....	245
14.3.8	re.split(pattern, string[, maxsplit=0, flags=0]).....	246
14.3.9	re.sub(pattern, repl, string[, count, flags]) .....	247
14.4	常用操作示例.....	248
14.4.1	验证.....	248
14.4.2	提取.....	248
14.4.3	替换.....	249
14.4.4	切分.....	250
第 15 章	Ruby.....	251
15.1	预备知识.....	251
15.2	正则功能详解.....	252
15.2.1	列表.....	252

15.2.2	字符组.....	252
15.2.3	Unicode 属性 .....	253
15.2.4	字符组简记法.....	254
15.2.5	单词边界.....	254
15.2.6	行起始/结束位置.....	255
15.2.7	环视.....	256
15.2.8	匹配模式.....	256
15.2.9	捕获分组的引用.....	257
15.3	正则 API 简介 .....	257
15.3.1	Regexp .....	257
15.3.2	Regexp.match(text) .....	259
15.3.3	Regexp.quote(text)和 Regexp.escape(text).....	260
15.3.4	String.index(Regexp) .....	261
15.3.5	String.scan(Regexp) .....	261
15.3.6	String.slice(Regexp).....	262
15.3.7	String.split(Regexp) .....	262
15.3.8	String.sub(Regexp, Str) .....	263
15.3.9	String.gsub(Regexp, String) .....	264
15.4	常用操作示例.....	264
15.4.1	验证.....	264
15.4.2	提取.....	265
15.4.3	替换.....	265
15.4.4	切分.....	265
15.5	Ruby 1.9 的新变化 .....	266
第 16 章	Linux/UNIX.....	268
16.1	POSIX.....	268
16.1.1	POSIX 规范 .....	268
16.1.2	POSIX 字符组 .....	269
16.2	vi .....	271
16.2.1	字符组及简记法.....	271
16.2.2	量词.....	272
16.2.3	多选结构和捕获分组.....	272
16.2.4	环视.....	273
16.2.5	锚点和单词边界.....	273
16.2.6	替换操作的特殊字符 .....	274

16.2.7	replacement 中的特殊变量 .....	276
16.2.8	补充.....	276
16.3	grep .....	277
16.3.1	基本用法.....	277
16.3.2	字符组.....	277
16.3.3	锚点和单词边界.....	278
16.3.4	量词.....	278
16.3.5	多选结构和捕获分组 .....	279
16.3.6	options.....	279
16.3.7	egrep 和 fgrep .....	280
16.3.8	补充.....	280
16.4	awk.....	281
16.4.1	基本用法.....	281
16.4.2	字符组及简记法.....	282
16.4.3	锚点和单词边界.....	283
16.4.4	量词.....	283
16.4.5	多选结构.....	284
16.4.6	补充.....	284
16.5	sed .....	284
16.5.1	基本用法.....	284
16.5.2	字符组及简记法.....	285
16.5.3	锚点和单词边界.....	285
16.5.4	量词.....	286
16.5.5	多选结构和捕获分组 .....	286
16.5.6	options.....	286
16.5.7	补充.....	287
16.6	总结 .....	288
附录 A	常用语言中正则特性一览 .....	291
附录 B	常用的正则表达式 .....	293
附录 C	常用的正则表达式工具及资源 .....	309

# 第一部分

# 第1章 字符组

## 1.1 普通字符组

**字符组**（Character Class）<sup>1</sup>是正则表达式最基本的结构之一，要理解正则表达式的“灵活”，认识它是第一步。

顾名思义，字符组就是一组字符，在正则表达式中，它表示“在同一个位置可能出现的各种字符”，其写法是在一对方括号[和]之间列出所有可能出现的字符，简单的字符组比如[ab]、[314]、[#.?]在解决一些常见问题时，使用字符组可以大大简化操作，下面举“匹配数字字符”的例子来说明。

字符可以分为很多类，比如数字、字母、标点等。有时候要求“只出现一个数字字符”，换句话说，这个位置上的字符只能是0、1、2、…、8、9这10个字符之一。要进行这种判断，通常的思路是：用10个条件分别判断字符是否等于这10个字符，对10个结果取“或”，只要其中一个条件成立，就返回True，表示这是一个数字字符，其伪代码如下例1-1所示。

### 例 1-1 判断数字字符的伪代码

```
charStr == "0" || charStr == "1" ... || charStr == "9"
```

注：因为正则表达式处理的都是“字符串”（String）而不是“字符”，所以这里假设变量charStr（虽然它只包含一个字符）也是字符串类型，使用了双引号，在有些语言中字符串也用单引号表示。

这种解法的问题在于太烦琐——如果要判断是否是一个小写英文字母，就要用||连接26个判断；如果还要兼容大写字母，则要连接52个判断，代码长到几乎无法阅读。相反，用字符组解决起来却异常简单，具体思路是：列出可能出现的所有字符（在这个例子里就是10个数字字符），只要出现了其中任何一个，就返回True。例1-2给出了使用字符组判断的例子，程序语言使用

---

<sup>1</sup> 在有的资料中，写作Character Set，所以也有人翻译为“字符类”或者“字符集”。不过在计算机术语中，“类”是和“对象”相关的，“字符集”常常表示Character Set（比如GBK、UTF-8之类），所以本书中没有采用这两个名字。

Python。

### 例 1-2 用正则表达式判断数字字符

```
re.search("[0123456789]", charStr) != None
```

`re.search()` 是 Python 提供的正则表达式操作函数，表示“进行正则表达式匹配”；`charStr` 仍然是需要判断的字符串，而 `[0123456789]` 则是以字符串形式给出的正则表达式，它是一个字符组，表示“这里可以是 0、1、2、…、8、9 中的任意一个字符。只要 `charStr` 与其中任何一个字符相同（或者说“`charStr` 可以由 `[0123456789]` 匹配”），就会得到一个 `MatchObject` 对象（这个对象暂时不必关心，在第 21 页会详细讲解）；否则，返回 `None`。所以判断结果是否为 `None`，就可以判断 `charStr` 是否是数字字符。

当今流行的编程语言大多支持正则表达式，上面的例子在各种语言中的写法大抵相同，唯一的区别在于如何调用正则表达式的功能，所以用法其实大同小异。例 1-3 列出了常见语言中的表示，如果你现在就希望知道语言的细节，可以参考本书第三部分的具体章节。

### 例 1-3 用正则表达式判断数字字符在各种语言中的应用<sup>1</sup>

#### .NET (C#)

```
//能匹配则返回 true，否则返回 false
Regex.IsMatch(charStr, "[0123456789]);
```

#### Java

```
//能匹配则返回 true，否则返回 false
charStr.matches("[0123456789]);
```

#### JavaScript

```
//能匹配则返回 true，否则返回 false
/[0123456789]/.test(charStr);
```

#### PHP

```
//能匹配则返回 1，否则返回 0
preg_match("/[0123456789]/", charStr);
```

#### Python

```
#能匹配则返回 RegexObject，否则返回 None
re.search("[0123456789]", charStr)
```

#### Ruby

```
#能匹配则返回 0，否则返回 nil
charStr =~ /[0123456789]/
```

---

<sup>1</sup> 客观地说，Perl 是正则表达式处理最方便的编程语言，考虑到今天使用 Perl 的人数，以及 Perl 程序员一般都熟练掌握正则表达式的现实，本书没有给出 Perl 语言的例子。



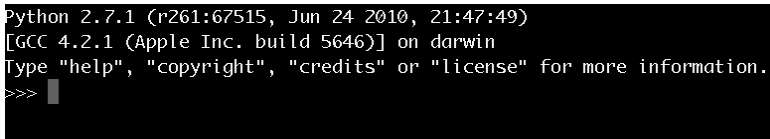
可以看到，不同语言使用正则表达式的方法也不相同。如果仔细观察会发现 Java、.NET、Python、PHP 中的正则表达式，都要以字符串形式给出，两端都有双引号<sup>1</sup>；而 Ruby 和 JavaScript 中的正则表达式则不必如此，只在首尾有两个斜线字符`/`，这也是不同语言中使用正则表达式的不同之处。不过，这个问题现在不需要太关心，因为本书中大部分例子以 Python 程序来讲解，下面讲解关于 Python 的基础知识，其他语言的细节留到后文会详细介绍。

## 1.2 关于 Python 的基础知识

本书选择使用 Python 语言来演示实际的匹配结果，因为它能在多种操作系统中运行，安装也很方便；另一方面，Python 是解释型语言，输入代码就能看到结果，方便动手实践。考虑到不是所有人都熟悉 Python，这里专门用一节来介绍。

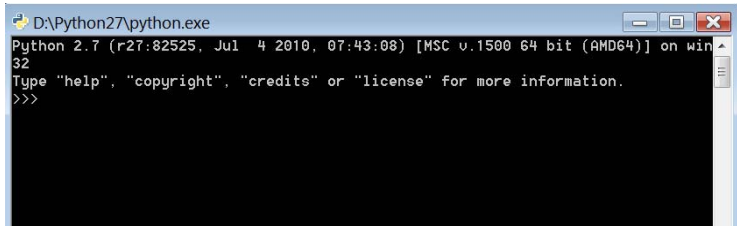
如果你的机器上没有安装 Python，可以从 <http://python.org/download/> 下载，目前 Python 有 2 和 3 两个版本，本书的例子以 2 版本为准<sup>1</sup>。请选择自己平台对应的程序下载并安装（目前 MacOS、Linux 的各种发行版一般带有 Python，具体可以在命令行下输入 `python`，看是否启动对应的程序）。

然后可以启动 Python，在 MacOS 和 Linux 下是输入 `python`，会显示出 Python 提示符，进入交互模式，如图 1-1（Linux 下的提示符与 MacOS 下的差不多，所以此处不列出）；而在 Windows 下，需要在“开始”菜单的“程序”中，选择 Python 目录下的 Python(command line)，如图 1-2 所示。



```
Python 2.7.1 (r261:67515, Jun 24 2010, 21:47:49)
[GCC 4.2.1 (Apple Inc. build 5646)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

图 1-1 MacOS 下的 Python 提示符



```
D:\Python27\python.exe
Python 2.7 (r27:82525, Jul 4 2010, 07:43:08) [MSC v.1500 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

图 1-2 Windows 下的 Python 提示符

---

<sup>1</sup> 本书写作时，2.x 最新的版本为 2.7，示例以此为准。Python 3 虽然已经正式发行，但相对 2.x 变化较大，而 2.x 较为流行，所以采用了 2.x 版本。关于 2.x 和 3.x 的差别，在 Python 一章有详细介绍。

Python 中常用的关于正则表达式的函数是 `re.search()`，使用它必须首先导入正则表达式对应的包（package），也就是输入下面的代码。

```
#导入正则表达式对应的包
import re
```

通常的用法是提供两个参数：`re.search(pattern, string)`，其中 *pattern* 是字符串形式提供的正则表达式，*string* 是需要匹配的字符串；如果能匹配，则返回一个 `MatchObject`（详细介绍请参考第 241 页，暂时可以不必关心），这时提示符会显示类似 `<_sre.SRE_Match object at 0x000000001D8E578>` 之类的结果；如果不能匹配，结果是 `None`（这是 Python 中的一个特殊值，类似其他某些语言中的 `Null`），不会有任何显示。图 1-3 演示了运行 Python 语句的结果。

```
>>> import re
>>> re.search("[0123456789]", "4")
<_sre.SRE_Match object at 0x000000001D8E578>
>>> re.search("[0123456789]", "a")
>>>
```

图 1-3 观察 `re.search()` 匹配的返回值

注：`>>>` 是等待输入的提示符，以 `>>>` 开头的行，之后文本是用户输入的语句；其他行是系统生成的，比如打印出语句的结果（在交互模式下，匹配结果会自动输出，便于观察；真正程序运行时不会如此）。

为讲解清楚、形象、方便，本书中的程序部分需要做两点修改。

第一，因为暂时还不需要关心匹配结果的细节，只关心有没有结果，所以在 `re.search()` 之后添加判断返回值是否为 `None`，如果为 `True`，则表示匹配成功，否则返回 `False` 表示匹配失败。为节省版面，尽可能用注释表示这类匹配结果，如 `# ==> True` 或者 `# ==> False`，附在语句之后。

第二，目前我们关心的是整个字符串是否能由正则表达式匹配。但是，在默认情况下 `re.search(pattern, string)` 只判断 *string* 的某个子串能否由 *pattern* 匹配，即便 *pattern* 只能匹配 *string* 的一部分，也不会返回 `None`。为了测试整个 *string* 能否由 *pattern* 匹配，在 *pattern* 两端加上 `^` 和 `$`。`^` 和 `$` 是正则表达式中的特殊字符，它们并不匹配任何字符，只是表示“定位到字符串的起始位置”和“定位到字符串的结束位置”（原理如图 1-4 所示，如果你现在就希望详细了解这两个特殊字符，可以参考第 61 页），这样就保证：只有在整个 *string* 都可以由 *pattern* 匹配时，才算匹配成功，不返回 `None`，如例 1-4 所示。

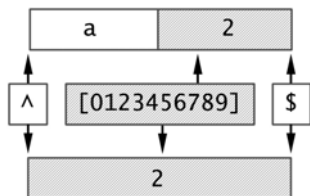


图 1-4 `^[0123456789]$` 的匹配

例 1-4 使用`^`和`$`测试 string 由 pattern 完整匹配

```
# 只要字符串中包含数字字符, 就可以匹配
re.search("[0123456789]", "2") != None      # => True
re.search("^[0123456789]$", "12") != None   # => False
re.search("[0123456789]", "a2") != None      # => True

# 整个字符串就是一个数字字符, 才可以匹配
re.search("[0123456789]", "2") != None      # => True
re.search("^[0123456789]$", "12") != None    # => False
re.search("^[0123456789]$", "a2") != None    # => False
```

## 1.3 普通字符组 (续)

介绍完关于 Python 的基础知识, 继续讲解字符组。字符组中的字符排列顺序并不影响字符组的功能, 出现重复字符也不会影响, 所以 `[0123456789]` 完全等价于 `[9876543210]`、`[1029384756]`、`[9988876543210]`。

不过, 代码总是要容易编写, 方便阅读, 正则表达式也是一样, 所以一般并不推荐在字符组中出现重复字符。而且, 还应该让字符组中的字符排列更符合认知习惯, 比如 `[0123456789]` 就好过 `[0192837465]`。为此, 正则表达式提供了 **-范围表示法** (range), 它更直观, 能进一步简化字符组。

所谓“-范围表示法”, 就是用 `[x-y]` 的形式表示 `x` 到 `y` 整个范围内的字符, 省去一一列出的麻烦, 这样 `[0123456789]` 就可以表示为 `[0-9]`。如果你觉得这不算什么, 那么确实比 `[abcdefghijklmnopqrstuvwxyz]` 简单太多了。

你可能会问,“-范围表示法”的范围是如何确定的? 为什么要写作 `[0-9]`, 而不写作 `[9-0]`?

要回答这个问题, 必须了解范围表示法的实质。在字符组中, `-` 表示的范围, 一般是根据字符对应的 **码值** (Code Point, 也就是字符在对应编码表中的编码的数值) 来确定的, 码值小的字符在前, 码值大的字符在后。在 ASCII 编码中 (包括各种兼容 ASCII 的编码中), 字符 `0` 的码值是 48 (十进制), 字符 `9` 的码值是 57 (十进制), 所以 `[0-9]` 等价于 `[0123456789]`; 而 `[9-0]` 则是错误的范围, 因为 `9` 的码值大于 `0`, 所以会报错。程序代码见例 1-5。

例 1-5 `[0-9]`是合法的, `[9-0]`会报错

```
re.search("^[0-9]$", "2") != None          # => True

re.search("^[9-0]$", "2") != None
Traceback (most recent call last):
error: bad character range
```

如果知道 `0~9` 的码值是 48~57, `a~z` 的码值是 97~122, `A~Z` 的码值是 65~90, 能不能用 `[0-z]` 统一表示数字字符、小写字母、大写字母呢？

答案是：勉强可以，但不推荐这么做。根据惯例，字符组的范围表示法都表示一类字符（数字字符是一类，字母字符也是一类），所以虽然 `[0-9]`、`[a-z]` 都是很好理解的，但 `[0-z]` 却很难理解，不熟悉 ASCII 编码表的人甚至不知道这个字符组还能匹配大写字母，更何况，在码值 48 到 122 之间，除去数字字符（码值 48~57）、小写字母（码值 97~122）、大写字母（码值 65~90），还有不少标点符号（参见表 1-1），从字符组 `[0-z]` 中却很难看出来，使用时就容易引起误会，例 1-6 所示的程序就很可能让人莫名其妙。

表 1-1 ASCII 编码表（片段）

码值	字符	码值	字符	码值	字符	码值	字符	码值	字符
48	0	63	?	78	N	93	]	108	l
49	1	64	@	79	O	94	^	109	m
50	2	65	A	80	P	95	_	110	n
51	3	66	B	81	Q	96	`	111	o
52	4	67	C	82	R	97	a	112	p
53	5	68	D	83	S	98	b	113	q
54	6	69	E	84	T	99	c	114	r
55	7	70	F	85	U	100	d	115	s
56	8	71	G	86	V	101	e	116	t
57	9	72	H	87	W	102	f	117	u
58	:	73	I	88	X	103	g	118	v
59	;	74	J	89	Y	104	h	119	w
60	<	75	K	90	Z	105	i	120	x
61	=	76	L	91	[	106	j	121	y
62	>	77	M	92	\	107	k	122	z

例 1-6 `[0-z]` 的奇怪匹配

```
re.search("^[0-z]$", "A") != None      # => True
re.search("^[0-z]$", ":") != None      # => True
```

在字符组中可以同时并列多个“-范围表示法”，字符组 `[0-9a-zA-Z]` 可以匹配数字、大写字母或小写字母；字符组 `[0-9a-zA-Z]` 可以匹配数字、大、小写形式的 `a~f`，它可以用来验证十六进制字符，代码见例 1-7。

## 例 1-7 [0-9a-fA-F]准确判断十六进制字符

```
re.search("[0-9a-fA-F]", "0") != None # => True
re.search("[0-9a-fA-F]", "c") != None # => True
re.search("[0-9a-fA-F]", "i") != None # => False
re.search("[0-9a-fA-F]", "C") != None # => True
re.search("[0-9a-fA-F]", "G") != None # => False
```

在不少语言中，还可以用转义序列 `\xhex` 来表示一个字符，其中 `\x` 是固定前缀，表示转义序列的开头，`num` 是字符对应的码值（Code Point，详见第 127 页，下文用 `U+127` 表示），是一个两位的十六进制数值。比如字符 `A` 的码值是 41（十进制则为 65），所以也可以用 `\x41` 表示。

字符组中有时会出现这种表示法，它可以表现一些难以输入或者难以显示的字符，比如 `\x7F`；也可以用来方便地表示某个范围，比如所有 ASCII 字符对应的字符组就是 `[\x00-\x7F]`，代码见例 1-8。这种表示法很重要，在第 120 页还会讲到它，依靠这种表示法可以很方便地匹配所有的中文字符。

## 例 1-8 [\x00-\x7F]准确判断 ASCII 字符

```
re.search("[\x00-\x7F]", "c") != None # => True
re.search("[\x00-\x7F]", "I") != None # => True
re.search("[\x00-\x7F]", "0") != None # => True
re.search("[\x00-\x7F]", "<") != None # => True
```

## 1.4 元字符与转义

在上面的例子里，字符组中的横线 `-` 并不能匹配横线字符，而是用来表示范围，这类字符叫做元字符（meta-character）。字符组的开方括号 `[`、闭方括号 `]` 和之前出现的 `^`、`$` 都算元字符。在匹配中，它们有着特殊的意义。但是，有时候并不需要表示这些特殊意义，只需要表示普通字符（比如“我就想表示横线字符 `-`”），此时就必须做特殊处理。

先来看字符组中的 `-`，如果它紧邻着字符组中的开方括号 `[`，那么它就是普通字符，其他情况下都是元字符；而对于其他元字符，取消特殊含义的做法都是转义，也就是在正则表达式中的元字符之前加上反斜线字符 `\`。

如果要在字符组内部使用横线 `-`，最好的办法是将它排列在字符组的最开头。`[-09]` 就是包含三个字符 `-`、`0`、`9` 的字符组；`[0-9]` 是包含 `0~9` 这 10 个字符的字符组，`[-0-9]` 则是由“范围表示法” `0-9` 和横线 `-` 共同组成的字符组，它可以匹配 11 个字符，例 1-9 说明了使用横线 `-` 的各种情况。

### 例 1-9 -出现在不同位置，含义不同

```
#作为普通字符
re.search("^[-09]$", "3") != None      # => False
re.search("^[-09]$", "-") != None      # => True
#作为元字符
re.search("^[-0-9]$", "3") != None      # => True
re.search("^[-0-9]$", "-") != None      # => False
#转义之后作为普通字符
re.search("^[\-0-9]$", "3") != None     # => False
re.search("^[\-0-9]$", "-") != None     # => True
```

仔细观察会发现，在正文里说“在正则表达式中的元字符之前加上反斜线字符`\`”，而在代码里写的却不是`[0-9]`，而是`[0\-9]`。这并不是输入错误。

因为在这段程序里，正则表达式是以字符串（String）的方式<sup>1</sup>提供的，而字符串本身也有关于转义的规定（你或许记得，在字符串中有`\n`、`\t`之类的转义序列）。上面说的“正则表达式”，其实是经过“字符串转义处理”之后的字符串的值，正则表达式`[0-9]`包含 6 个字符：`[`、`0`、`\`、`-`、`9`、`]`，在字符串中表达这 6 个字符；但是在源代码里，必须使用 7 个字符：`\`需要转义成`\\`，因为处理字符串时，反斜线和之后的字符会被认为是转义序列（Escape Sequence），比如`\n`、`\t`都是合法的转义序列，然而`\-`不是。

这个问题确实有点麻烦。正则表达式是用来处理字符串的，但它又不完全等于字符串，正则表达式中的每个反斜线字符`\`，在字符串中（也就是正则表达式之外）还必须转义为`\\`。所以之前所说的是“正则表达式`[0-9]`”，程序里写的却是`[0\\-9]`，这确实有点麻烦。

不过，Python 提供了**原生字符串**（Raw String），它非常适合于正则表达式：正则表达式是怎样，原生字符串就是怎样，完全不需要考虑正则表达式之外的转义（只有双引号字符是例外，原生字符串内的双引号字符必须转义写成`\"`）。原生字符串的形式是 `r"string"`，也就是在普通字符串之前添加 `r`，示例代码如例 1-10。

### 例 1-10 原生字符串的使用

```
#原生字符串和字符串的等价
r"^[-0-9]$" == "^[\-0-9]$"             # => True
#原生字符串的转义要简单许多
re.search(r"^[-0-9]$", "3") != None     # => False
re.search(r"^[-0-9]$", "-") != None     # => True
```

原生字符串清晰易懂，省去了烦琐的转义，所以从现在开始，本书中的 Python 示范代码都

---

<sup>1</sup> 具体来说，在 Java、PHP、Python、.NET 等语言中，正则表达式都是以字符串的形式给出的，在 Ruby 和 JavaScript 中则不是这样。详细的说明，请参考第 96 页。

会使用原生字符串来表示正则表达式。另外，.NET 和 Ruby 中也有原生字符串，也有一些语言并没有提供原生字符串（比如 Java），所以在第 6 章（☞93）会专门讲解转义问题。不过，现在只需要知道 Python 示范代码中使用了原生字符串即可。

继续看转义，如果希望在字符组中列出闭方括号`]`，比如`[012]345]`，就必须在它之前使用反斜线转义，写成`[012\]345]`；否则，结果就如例 1-11 所示，正则表达式将`]`与最近的`[`匹配，这个表达式就成了“字符组`[012]`加上 4 个字符`345]`”，它能匹配的是字符串`0345]`或`1345]`或`2345]`，却不能匹配`]`。

例 1-11 `]`出现在不同位置，含义不同

```
#未转义的]
re.search(r"^[012]345]$", "2345") != None      # => True
re.search(r"^[012]345]$", "5") != None          # => False
re.search(r"^[012]345]$", "]") != None          # => False
#转义的]
re.search(r"^[012\]345]$", "2345") != None      # => False
re.search(r"^[012\]345]$", "5") != None          # => True
re.search(r"^[012\]345]$", "]") != None          # => True
```

除去字符组内部的`-`，其他元字符的转义都必须在字符之前添加反斜线，`[`的转义也是如此。如果只希望匹配字符串`012]`，直接使用正则表达式`[012]`是不行的，因为这会被识别为一个字符组，它只能匹配`0`、`1`、`2`这三个字符中的任意一个；而必须转义，把正则表达式写作`\[012]`，请注意，只有开方括号`[`需要转义，闭方括号`]`不需要转义，如例 1-12 所示。

例 1-12 取消其他元字符的特殊含义

```
re.search(r"^[012]345]$", "3") != None          # => False
re.search(r"^[012\]345]$", "3") != None          # => True
re.search(r"^[012]$", "[012]") != None          # => False
re.search(r"^[012\]$", "[012]") != None          # => True
```

## 1.5 排除型字符组

在方括号`[...]`中列出希望匹配的所有字符，这种字符组叫做“普通字符组”，它的确非常方便。不过，也有些问题是普通字符组不能解决的。

给定一个由两个字符构成的字符串 `str`，要判断这两个字符是否都是数字字符，可以用`[0-9][0-9]`来匹配。但是，如果要求判断的是这样的字符串——第一个字符不是数字字符，第二个字符才是数字字符（比如 `A8`、`x6`）<sup>1</sup>——应当如何办？数字字符的匹配很好处理，用`[0-9]`

<sup>1</sup> 一般来说，计算机中的偏移值都是从 0 开始的。此处考虑到叙述自然，使用了“第一个字符”和“第二个字符”的说法，其中“第一个字符”指最左端，也就是偏移值为 0 的字符；“第二个字符”指紧跟在它右侧，也就是偏移值为 1 的字符。

即可；“不是数字”则很难办——不是数字的字符太多了，全部列出几乎不可能，这时就应当使用排除型字符组。

**排除型字符组**（Negated Character Class）非常类似普通字符组`[...]`，只是在开方括号`[`之后紧跟一个脱字符`^`，写作`[^...]`，表示“在当前位置，匹配一个没有列出的字符”。所以`[^0-9]`就表示“`0~9`之外的字符”，也就是“非数字字符”。那么，`[^0-9][0-9]`就可以解决问题了，如例 1-13 所示。

#### 例 1-13 使用排除型字符组

```
re.search(r"^[^0-9][0-9]$", "A8") != None      # => True
re.search(r"^[^0-9][0-9]$", "x6") != None      # => True
```

排除型字符组看起来很简单，不过新手常常会犯一个错误，就是把“在当前位置，匹配一个没有列出的字符”理解成“在当前位置不要匹配列出的字符”两者其实是不同的，后者暗示“这里不出现任何字符也可以”。例 1-14 很清楚地说明：**排除型字符组必须匹配一个字符**，这点一定要记住。

#### 例 1-14 排除型字符组必须匹配一个字符

```
re.search(r"^[^0-9][0-9]$", "8") != None      # => False
re.search(r"^[^0-9][0-9]$", "A8") != None     # => True
```

除了开方括号`[`之后的`^`，排除型字符组的用法与普通字符组几乎完全相同，唯一需要改动的是：在排除型字符组中，如果需要表示横线字符`-`（而不是用于“-范围表示法”），那么`-`应该紧跟在`^`之后；而在普通字符组中，作为普通字符的横线`-`应该紧跟在开方括号之后，如例 1-15 所示。

#### 例 1-15 在排除型字符组中，紧跟在`^`之后的`-`不是元字符

```
#匹配一个-、0、9 之外的字符
re.search(r"^[^0-9]$", "-") != None           # => False
re.search(r"^[^0-9]$", "8") != None           # => True

#匹配一个 0~9 之外的字符
re.search(r"^[^0-9]$", "-") != None           # => True
re.search(r"^[^0-9]$", "8") != None           # => False
```

在排除型字符组中，`^`是一个元字符，但只有它紧跟在`[`之后时才是元字符，如果想表示“这个字符组中可以出现`^`字符”，不要让它紧挨着`[`即可，否则就要转义。例 1-16 给出了三个正则表达式，后两个表达式实质是一样的，但第三种写法很麻烦，理解起来也麻烦，不推荐使用。



## 例 1-16 排除型字符组的转义

```
#匹配一个 0、1、2 之外的字符
re.search(r"^[^012]$", "^") != None      # => True
#匹配 4 个字符之一: 0、^、1、2
re.search(r"^[0^12]$", "^") != None      # => True
#^紧跟在[]之后, 但经过转义变为普通字符, 等于上一个表达式, 不推荐
re.search(r"^[^012]$", "^") != None      # => True
```

## 1.6 字符组简记法

用 `[0-9]`、`[a-z]` 等字符组, 可以很方便地表示数字字符和小写字母字符。对于这类常用的字符组, 正则表达式提供了更简单的记法, 这就是**字符组简记法** (shorthands)。

常见的字符组简记法有 `\d`、`\w`、`\s`。从表面上看, 它们与 `[...]` 完全没联系, 其实是一致的。其中 `\d` 等价于 `[0-9]`, 其中的 d 代表“数字 (digit)”; `\w` 等价于 `[0-9a-zA-Z_]`, 其中的 w 代表“单词字符 (word)”; `\s` 等价于 `[\t\r\n\v\f]` (第一个字符是空格), s 表示“空白字符 (space)”。例 1-17 说明了这几个字符组简记法的典型匹配。

例 1-17 字符组简记法 `\d`、`\w`、`\s`

```
#如果没有原生字符串, \d 就必须写作\\d
re.search(r"^\d$", "8") != None          # => True
re.search(r"^\d$", "a") != None          # => False

re.search(r"^\w$", "8") != None          # => True
re.search(r"^\w$", "a") != None          # => True
re.search(r"^\w$", "_") != None          # => True

re.search(r"^\s$", " ") != None          # => True
re.search(r"^\s$", "\t") != None         # => True
re.search(r"^\s$", "\n") != None         # => True
```

一般印象中, 单词字符似乎只包含大小写字母, 但是字符组简记法中的“单词字符”不只有大小写单词, 还包括数字字符和下划线 `_`, 其中的下划线 `_` 尤其值得注意: 在进行数据验证时, 有可能只容许输入“数字和字母”, 有人会偷懒用 `\w` 验证, 而忽略了 `\w` 能匹配下划线, 所以这种匹配并不严格, `[0-9a-zA-Z]` 才是准确的选择。

“空白字符”并不难定义, 它可以是空格字符、制表符 `\t`, 回车符 `\r`, 换行符 `\n` 等各种“空白”字符, 只是不方便展现 (因为显示和印刷出来都是空白)。不过这也提醒我们注意, 匹配时看到的“空白”可能不是空格字符, 因此, `\s` 才是准确的选择。

字符组简记法可以单独出现, 也可以使用在字符组中, 比如 `[0-9a-zA-Z]` 也可以写作 `[\da-zA-Z]`, 所以匹配十六进制字符的字符组可以写成 `[\da-fA-F]`。字符组简记法也可以用

在排除型字符组中，比如`[^0-9]`就可以写成`[^\d]`，`[^0-9a-zA-Z_]`就可以写成`[^\w]`，代码如例 1-18。

### 例 1-18 字符组简记法与普通字符组混用

```
#用在普通字符组内部
re.search(r"^[da-zA-Z]$", "8") != None    # => True
re.search(r"^[da-zA-Z]$", "a") != None    # => True
re.search(r"^[da-zA-Z]$", "C") != None    # => True
#用在排除型字符组内部
re.search(r"^[^\w]$", "8") != None        # => False
re.search(r"^[^\w]$", "_") != None        # => False
re.search(r"^[^\w]$", ",") != None        # => True
```

相对于`\d`、`\w`和`\s`这三个普通字符组简记法，正则表达式也提供了对应排除型字符组的简记法：`\D`、`\W`和`\S`——字母完全一样，只是改为大写。这些简记法匹配的字符互补：`\s`能匹配的字符，`\S`一定不能匹配；`\w`能匹配的字符，`\W`一定不能匹配；`\d`能匹配的字符，`\D`一定不能匹配。例 1-19 示范了这几个字符组简记法的应用。

### 例 1-19 \D、\W、\S 的使用

```
#\d 和 \D
re.search(r"^\d$", "8") != None           # => True
re.search(r"^\d$", "a") != None           # => False
re.search(r"^\D$", "8") != None           # => False
re.search(r"^\D$", "a") != None           # => True
#\w 和 \W
re.search(r"^\w$", "c") != None           # => True
re.search(r"^\w$", "!") != None           # => False
re.search(r"^\W$", "c") != None           # => False
re.search(r"^\W$", "!") != None           # => True
#\s 和 \S
re.search(r"^\s$", "\t") != None          # => True
re.search(r"^\s$", "0") != None           # => False
re.search(r"^\S$", "\t") != None          # => False
re.search(r"^\S$", "0") != None           # => True
```

妥善利用这种互补的属性，可以得到一些非常巧妙的效果，最简单的应用就是字符组`[\s\S]`。初看起来，在同一个字符组中并列两个互补的简记法，这种做法有点奇怪，不过仔细想想就会明白，`\s`和`\S`组合在一起，匹配的就是“所有的字符”（或者叫“任意字符”）。许多语言中的正则表达式并没有直接提供“任意字符”的表示法，所以`[\s\S]`、`[\w\W]`、`[\d\D]`虽然看起来有点古怪，但确实可以匹配任意字符<sup>1</sup>。

<sup>1</sup> 许多关于正则表达式的文档说：点号`.`能匹配“任意字符”。但在默认情况下，点号其实不能匹配换行符，具体请参考第 84 页。

关于字符组简记法，最后需要补充两点：第一，如果字符组中出现了字符组简记法，最好不要出现单独的`-`，否则可能引起错误，比如`[\d-a]`就很让人迷惑，在有些语言中，`-`会被作为普通字符，而在有些语言中，这样写会报错；第二，以上说的`\d`、`\w`、`\s`的匹配规则，都是针对ASCII编码而言的，也叫**ASCII匹配规则**。但是，目前一些语言中的正则表达式已经支持了Unicode字符，那么数字字符、单词字符、空白字符的范围，已经不仅仅限于ASCII编码中的字符。关于这个问题，具体细节在后文有详细的介绍，如果你现在就想知道，可以翻到第115页。

## 1.7 字符组运算

以上介绍了字符组的基本功能，它们在常用的语言中都有提供；还有些语言中为字符组提供了更强大的功能，比如Java和.NET就提供了字符组运算的功能，可以在字符组内进行集合运算，在某些情况下这种功能非常实用。

如果要匹配所有的元音字母（为讲解简单考虑，暂时只考虑小写字母的情况），可以用`[aeiou]`，但是要匹配所有的辅音字母却没有那么方便的办法，最直接的写法是`[b-df-hj-np-tv-z]`，不但烦琐，而且难理解。其实，从26个字母中“减去”元音字母，剩下的就是辅音字母，如果有办法做这个“减法”，就方便多了。

Java语言中提供了这样的字符组：`[[a-z]&&[^aeiou]]`，虽然初看有点古怪，但仔细看看，也不难理解。`[a-z]`表示26个英文字母，`[^aeiou]`表示除元音字母之外的所有字符（还包括大写字母、数字和各种符号），两者取交集，就得到“26个英文字母中，除去5个元音字母，剩下的21个辅音字母”。

.NET中也有这样的功能，只是写法不一样。同样是匹配辅音字母的字符组，.NET中写作`[a-z-[aeiou]]`，其逻辑是：从`[a-z]`能匹配的26个字符中，“减去”`[aeiou]`能匹配的元音字母。相对于Java，这种逻辑更符合直觉，但写法却有点古怪——不是`[[a-z]-[aeiou]]`，而是`[a-z-[aeiou]]`。例1-20集中演示了Java和.NET中的字符组运算。

### 例 1-20 字符组运算

#### Java

```
"a".matches("^[[a-z]&&[^aeiou]]$"); // => True
"b".matches("^[[a-z]&&[^aeiou]]$"); // => False
```

#### .NET

```
Regex.IsMatch("^[a-z-[aeiou]]$", "a"); // => True
Regex.IsMatch("^[a-z-[aeiou]]$", "b"); // => False
```

# 1.8 POSIX 字符组

前面介绍了常用的字符组，但是在某些文档中，你可能会发现类似[:digit:]、[:lower:]之类的字符组，看起来不难理解（digit 就是“数字”，lower 就是“小写”），但又很奇怪，它们就是 **POSIX 字符组**（POSIX Character Class）。因为某些语言的文档中出现了这些字符组，为避免困惑，这里有必要做个简要介绍。如果只使用常用的编程语言，可以忽略文档中的 POSIX 字符组，也可以忽略本节；如果想了解 POSIX 字符组，或者需要在 Linux/UNIX 下的各种工具（sed、awk、grep 等）中使用正则表达式，最好阅读本节。

之前介绍的字符组，都属于 Perl 衍生出来的正则表达式流派（Flavor），这个流派叫做 PCRE（Per Compatible Regular Expression）。在此之外，正则表达式还有其他流派，比如 POSIX（Portable Operating System Interface for uniX），它是一系列规范，定义了 UNIX 操作系统应当支持的功能，其中也包括关于正则表达式的规范，[:digit:]之类的字符组就是遵循 POSIX 规范的字符组。

常见的[a-z]形式的字符组，在 POSIX 规范中仍然获得支持，它的准确名称是 **POSIX 方括号表达式**（POSIX bracket expression），主要用在 UNIX/Linux 系统中。POSIX 方括号表达式与之前所说的字符组最主要的差别在于：在 POSIX 字符组中，反斜线\不是用来转义的。所以 POSIX 方括号表达式[\d]只能匹配\和 d 两个字符，而不是[0-9]对应的数字字符。

为了解决字符组中特殊意义字符的转义问题，POSIX 方括号表达式规定：如果要在字符组中表达字符]（而不是作为字符组的结束标记），应当让它紧跟在字符组的开方括号之后，所以[la]能匹配的字符就是]或 a；如果要在字符组中标识字符-（而不是“-范围表示法”），就必须将它放在字符组的闭方括号]之前，所以[a-]能匹配的字符就是 a 或-。

另一方面，POSIX 规范还定义了 POSIX 字符组（POSIX character class），它大致等于之前介绍的字符组简记法，都是使用类似[:digit:]、[:lower:]之类有明确意义的记号表示某类字符。

表 1-2 简要介绍了 POSIX 字符组，注意表格中与其对应的是 ASCII 字符组，也就是能匹配的 ASCII 字符（ASCII 编码表中码值在 0~127 之间的字符）。因为 POSIX 规范中有一个重要概念：locale（通常翻译为“语言环境”），它是一组与语言和文化相关的设定，包括日期格式、货币币值、字符编码等。POSIX 字符组的意义会根据 locale 的变化而变化，表 1-2 介绍的只是这些 POSIX 字符组在 ASCII 编码中的意义；如果换用其他的 locale（比如使用 Unicode 字符集），它们的意义可能会发生变化，具体请参考第 129 页。

表 1-2 POSIX 字符组

POSIX 字符组	说明	ASCII 字符组	等价的 PCRE 简记法
[:alnum:]*	字母字符和数字字符	[a-zA-Z0-9]	
[:alpha:]	字母	[a-zA-Z]	

( 续表 )

POSIX 字符组	说明	ASCII 字符组	等价的 PCRE 简记法
<code>[ :ASCII: ]</code>	ASCII 字符	<code>[ \x00-\x7F ]</code>	
<code>[ :blank: ]</code>	空格字符和制表符	<code>[ \t ]</code>	
<code>[ :cntrl: ]</code>	控制字符	<code>[ \x00-\x1F\x7F ]</code>	
<code>[ :digit: ]</code>	数字字符	<code>[ 0-9 ]</code>	<code>\d</code>
<code>[ :graph: ]</code>	空白字符之外的字符	<code>[ \x21-\x7E ]</code>	
<code>[ :lower: ]</code>	小写字母字符	<code>[ a-z ]</code>	
<code>[ :print: ]</code>	类似 <code>[ :graph: ]</code> ，但包括空白字符	<code>[ \x20-\x7E ]</code>	
<code>[ :punct: ]</code>	标点符号	<code>[ [ ! " # \$ % &amp; ' ( ) * + , . / : ; &lt; = &gt; ? @ ^ _ ` {   } ~ - ]</code>	
<code>[ :space: ]</code>	空白字符	<code>[ \t\r\n\v\f ]</code>	<code>\s</code>
<code>[ :upper: ]</code>	大写字母字符	<code>[ A-Z ]</code>	
<code>[ :word: ]*</code>	字母字符	<code>[ A-Za-z0-9_ ]</code>	<code>\w</code>
<code>[ :xdigit: ]</code>	十六进制字符	<code>[ A-Fa-f0-9 ]</code>	

注：标记\*的字符组简记法并不是 POSIX 规范中的，但使用很多，一般语言中都提供，文档中也会出现。

POSIX 字符组的使用也与 PCRE 字符组简记法的使用有所不同，主要区别在于，PCRE 字符组简记法可以脱离方括号直接出现，而 POSIX 字符组必须出现在方括号内。所以同样是匹配数字字符，PCRE 中可以直接写`\d`，而 POSIX 字符组必须写成`[[:digit:]]`。

在本书介绍的 6 种语言中，Java、PHP、Ruby 支持使用 POSIX 字符组。

在 PHP 中可以直接使用 POSIX 字符组，但是 PHP 中的 POSIX 字符组只识别 ASCII 字符，也就是说，任何非 ASCII 字符（比如中文字符）都不能由任何一个 POSIX 字符组匹配。

Ruby 的情况稍微复杂一点。Ruby 1.8 中的 POSIX 字符组只能匹配 ASCII 字符，而且不支持`[ :word: ]`和`[ :ASCII: ]`；Ruby 1.9 中的 POSIX 字符组可以匹配 Unicode 字符，而且支持`[ :word: ]`和`[ :ASCII: ]`。

Java 中的情况更加复杂。POSIX 字符组`[[:name:]]`必须使用`\p{name}`的形式，其中 *name* 为 POSIX 字符组对应的名字，比如`[ :space: ]`就应当写作`\p{Space}`，请注意第一个字母要大写，其他 POSIX 字符组都是这样，只有`[ :xdigit: ]`要写作`\p{XDigit}`。并且 Java 中的 POSIX 字符组，只能匹配 ASCII 字符。

## 第2章 量词

### 2.1 一般形式

根据上一章的介绍，可以用字符组`[0-9]`或者`\d`匹配单个数字字符。现在用正则表达式来验证更复杂的字符串，比如大陆地区的邮政编码。

粗略来看，邮政编码并没有特殊的规定，只是6位数字构成的字符串，比如`201203`、`100858`，所以用正则表达式来表示就是`\d\d\d\d\d\d`，如例2-1所示，只有同时满足“长度是6个字符”和“每个字符都是数字”两个条件，匹配才成功（同样，这里不能忽略`^`和`$`）。

例2-1 匹配邮政编码

```
re.search(r"^\d\d\d\d\d\d$", "100859") != None      # => True
re.search(r"^\d\d\d\d\d\d$", "201203") != None      # => True

re.search(r"^\d\d\d\d\d\d$", "20A203") != None      # => False
re.search(r"^\d\d\d\d\d\d$", "20103") != None       # => False
re.search(r"^\d\d\d\d\d\d$", "2012036") != None     # => False
```

虽然这不难理解，但`\d`重复了6次，读写都不方便。为此，正则表达式提供了量词(quantifier)，比如上面匹配邮政编码的表达式，就可以如例2-2那样，简写为`\d{6}`，它使用阿拉伯数字，更简洁也更直观。

例2-2 使用量词简化字符组

```
re.search(r"^\d{6}$", "100859") != None            # => True
re.search(r"^\d{6}$", "201203") != None            # => True

re.search(r"^\d{6}$", "20A203") != None            # => False
re.search(r"^\d{6}$", "20103") != None             # => False
re.search(r"^\d{6}$", "2012036") != None           # => False
```

量词还可以表示不确定的长度，其通用形式是`{m,n}`，其中`m`和`n`是两个数字（有些人习惯在代码中的逗号之后添加空格，这样更好看，但是量词中的逗号之后绝不能有空格），它限定

之前的元素<sup>1</sup>能够出现的次数，*m* 是下限，*n* 是上限（均为闭区间）。比如 `\d{4,6}`，就表示这个数字字符串的长度最短是 4 个字符（“单个数字字符”至少出现 4 次），最长是 6 个字符。

如果不确定长度的上限，也可以省略，只指定下限，写成 `\d{m,}`，比如 `\d{4,}` 表示“数字字符串的长度必须在 4 个字符以上”。

量词限定的出现次数一般都有明确下限，如果没有，则默认为 0。有一些语言（比如 Ruby）支持 `{,n}` 的记法，这时候并不是“不确定长度的下限”，而是省略了“下限为 0”的情况，比如 `\d{,6}` 表示“数字字符串最多可以有 6 个字符”。不过，这种用法并不是所有语言中都通用的，比如 Java 就不支持这种写法，所以必须写明 `{0,n}`。我推荐的做法是：最好使用 `{0,n}` 的记法，因为它是广泛支持的。表 2-1 集中说明了这几种形式的量词，例 2-3 展示了它们的使用。

表 2-1 量词的一般形式

量词	说明
<code>{n}</code>	之前的元素必须出现 <i>n</i> 次
<code>{m,n}</code>	之前的元素最少出现 <i>m</i> 次，最多出现 <i>n</i> 次
<code>{m,}</code>	之前的元素最少出现 <i>m</i> 次，出现次数无上限
<code>{0,n}</code>	之前的元素可以不出现，也可以出现，最多出现 <i>n</i> 次（在某些语言中可以写为 <code>{,n}</code> ）

例 2-3 表示不确定长度的量词

```
re.search(r"^\d{4,6}$", "123") != None      # => False
re.search(r"^\d{4,6}$", "1234") != None     # => True
re.search(r"^\d{4,6}$", "123456") != None   # => True
re.search(r"^\d{4,6}$", "1234567") != None  # => False

re.search(r"^\d{4,}$", "123") != None        # => False
re.search(r"^\d{4,}$", "1234") != None      # => True
re.search(r"^\d{4,}$", "123456") != None    # => True

re.search(r"^\d{0,6}$", "12345") != None     # => True
re.search(r"^\d{0,6}$", "123456") != None   # => True
re.search(r"^\d{0,6}$", "1234567") != None  # => False
```

2.2 常用量词

`{m,n}` 是通用形式的量词，正则表达式还有三个常用量词，分别是 `+`、`?`、`*`。它们的形态虽然不同于 `{m,n}`，功能却是相同的（也可以把它们理解为“量词简记法”），具体说明见表 2-2。

<sup>1</sup> 在上一章提到，字符组是正则表达式的基本“结构”之一，而此处提到之前的“元素”，在此做一点解释。在本书中，“结构”一般指的是正则表达式所提供功能的记法。比如字符组就是一种结构，下一章要提到的括号也是一种结构；而“元素”指的是具体的正则表达式中的某个部分，比如某个具体表达式中的字符组 `[a-z]`，可以算作一个元素，“元素”也叫“子表达式”（sub-expression）。

表 2-2 常用量词

常用量词	{m,n}等价形式	说明
*	{0,}	可能出现，也可能不出现，出现次数没有上限
+	{1,}	至少出现 1 次，出现次数没有上限
?	{0,1}	至多出现 1 次，也可能不出现

在实际应用中，在很多情况只需要表示这三种意思，所以常用量词的使用频率要高于  $\{m,n\}$ ，下面分别说明。

大家都知道，美国英语和英国英语有些词的写法是不一样的，比如 `traveler` 和 `traveller`，如果希望“通吃” `traveler` 和 `traveller`，就要求第 2 个 `l` 是“至多出现 1 次，也可能不出现”的，正好使用 `?` 量词：`travell?er`，如例 2-4 所示。

例 2-4 量词?的应用

```
re.search(r"^travell?er$", "traveler") != None    # => True
re.search(r"^travell?er$", "traveller") != None  # => True
```

其实这样的情况还有很多，比如 `favor` 和 `favour`、`color` 和 `colour`。此外还有很多其他应用场合，比如 `http` 和 `https`，虽然是两个概念，但都是协议名，可以用 `https?` 匹配；再比如表示价格的字符串，有可能是 `100` 也有可能是 `¥100`，可以用 `¥?100` 匹配<sup>1</sup>。

量词也广泛应用于解析 HTML 代码。HTML 是一种“标签语言”，它包含各种各样的 `tag`（标签），比如 `<head>`、`<img>`、`<table>` 等，这些 `tag` 的名字各异，形式却相同：从 `<` 开始，到 `>` 结束，在 `<` 和 `>` 之间有若干字符，“若干”的意思是长度不确定，但不能为 0（`<>` 并不是合法的 `tag`），也不能是 `>` 字符<sup>2</sup>。如果要用一个正则表达式匹配所有的 `tag`，需要用 `<` 匹配开头的 `<`，用 `>` 匹配结尾的 `>`，用 `[^>]+` 匹配中间的“若干字符”，所以整个正则表达式就是 `<[>]+>`，程序如例 2-5 所示。

例 2-5 量词+的应用

```
re.search(r"^<[>]+>$", "<bld>") != None          # => True
re.search(r"^<[>]+>$", "</table>") != None        # => True
re.search(r"^<[>]+>$", "<>") != None              # => False
```

类似的，也可以使用正则表达式匹配双引号字符串。不同的是，双引号字符串的两个双引号

<sup>1</sup> 实际上，这个问题比较复杂，因为 `¥` 并不是一个 ASCII 字符，所以 `¥?` 可能会产生问题，具体情况请参考第 7 章。

<sup>2</sup> 如果你对 HTML 代码比较了解，可能会有疑问，假如 `tag` 内部出现 `>` 符号，怎么办？这种情况确实存在，比如 `<input name=txt value=">">`。以目前已经讲解的知识还无法解决这个问题，不过下一章就会给出它的解法。



之间可以没有任何字符，""也是一个完全合法的双引号字符串，应该使用量词\*，于是整个正则表达式就成了 "[^"]\*"，程序见例 2-6。

例 2-6 量词\*的应用

```
re.search(r"^\"[^\"]*$", "\"some\"") != None    # => True
re.search(r"^\"[^\"]*$", "\"\"") != None       # => True
```

注：字符串之中表示双引号需要转义写成\"，这并不是正则表达式中的规定，而是为字符串转义考虑。

量词的使用有很多学问，不妨多看几个 tag 匹配的例子：tag 可以粗略分为 open tag 和 close tag，比如<head>就是 open tag，而</html>就是 close tag；另外还有一类标签是 self-closing tag，比如<br/>。现在来看分别匹配这三类 tag 的正则表达式。

open tag 的特点是以<开头，然后是“若干字符”（但不能以/开头），最后是>，所以对应的正则表达式是<[^/][^>]\*>; 注意：因为[^/]必须匹配一个字符，所以“若干字符”中其他部分必须写成[^>]\*，否则它无法匹配名字为单个字符的标签，比如<b>。

close tag 的特点是以</开头，之后是/字符，然后是“若干字符（但不能以/开头）”，最后是>，所以对应的正则表达式是</[^>]+>;

self-closing tag 的特点是以<开头，中间是“若干字符”，最后是/>，所以对应的正则表达式是<[^>]+/>。注意：这里不是<[^>/]+/>，排除型字符组只排除>，而不排除/，因为要确认的只是在结尾的>之前出现/，如果写成<[^>/]+/>，则要求 tag 内部不能出现/，就无法匹配这类的 tag 了。

表 2-3 列出了匹配几类 tag 的表达式。

表 2-3 各类 tag 的匹配

匹配所有 tag 的表达式	tag 分类	匹配分类 tag 的表达式
<[^>]+>	open tag	<[^/][^>]*>
	close tag	</[^>]+>
	self-closing tag	<[^>/]+/>

对比表格中“匹配所有 tag 的表达式”和“匹配分类 tag 的表达式”，可以发现它们的模式是相近的，只是细节上有差异。也就是说，通过变换字符组和量词，可以准确控制正则表达式能匹配的字符串的范围，达到不同的目的。这其实是使用正则表达式时的一条根本规律：使用合适的结构（包括字符组和量词），精确表达自己的意图，界定能匹配的文本。

再仔细观察，你或许会发现，匹配 open tag 的表达式，也可以匹配 self-closing tag：<[^/][^>]\*>能够匹配<br/>，因为[^>]\*并不排除对/的匹配。那么将表达式改为

`<[^/][^>]*[^/]>`, 就保证匹配的 open tag 不会以`/>`结尾了。

不过这会产生新的问题：`<[^/][^>]*[^/]>`能匹配的 tag, 在`<`和`>`之间出现了两个`[^/]`, 上一章已经讲过, 排除型字符组表示“在当前位置, 匹配一个没有列出的字符”, 所以 tag 里的字符串必须至少包含两个字符, 这样就无法匹配`<u>`了。

仔细想想, 真正要表达的意思是, 在 tag 内部的字符串不能以`/`开头, 也不能以`/`结尾, 如果这个字符串只包含一个字符, 那么它既是开头, 又是结尾, 使用两个排除型字符组显然是不合适的, 看起来没办法解决了。实际上, 只是现有的知识还不够解决这个问题而已, 在第 68 页有这个问题的详细解法。

## 2.3 数据提取

正则表达式的功能很多, 除去之前介绍的验证 (字符串能否由正则表达式匹配), 还可以从某个字符串中提取出某个字符串能匹配的所有文本。

上一章提到, `re.search()`如果匹配成功, 返回一个 `MatchObject` 对象。这个对象包含了匹配的信息, 比如表达式匹配的结果, 可以像例 2-7 那样, 通过调用 `MatchObject.group(0)`来获得。这个方法以后详细介绍, 现在只需要了解一点: 调用它可以得到表达式匹配的文本。

### 例 2-7 通过 MatchObject 获得匹配的文本

#注意这里使用链式编程

```
print re.search(r"\d{6}", "ab123456cd").group(0)
123456

print re.search(r"^(<[^>]+>$", "<bold>").group(0)
<bold>
```

这里再介绍一个方法: `re.findall(pattern, string)`。其中 `pattern` 是正则表达式, `string` 是字符串。这个方法会返回一个数组, 其中的元素是在 `string` 中依次寻找 `pattern` 能匹配的文本。

以邮政编码的匹配为例, 假设某个字符串中包含两个邮政编码: `zipcode1:201203`, `zipcode2:100859`, 仍然使用之前匹配邮政编码的正则表达式`\d{6}`, 调用 `re.findall()`可以将这两个邮政编码提取出来, 如例 2-8。注意, 这次要去掉表达式首尾的`^`和`$`, 因为要使用正则表达式在字符串中寻找匹配, 而不是验证整个字符串能否由正则表达式匹配。

### 例 2-8 使用 re.findall()提取数据

```
print re.findall(r"\d{6}", "zipcode1:201203, zipcode2:100859")
['201203', '100859']
```

```
#也可以逐个输出
for zipcode in re.findall(r"\d{6}", "zipcode1:201203, zipcode2:100859"):
    print zipcode
201203
100859
```

借助之前的匹配各种 tag 的正则表达式, 还可以通过 `re.findall()` 将某个 HTML 页面中所有的 tag 提取出来, 下面以 Yahoo 首页为例。

首先要读入 `http://www.yahoo.com/` 的 HTML 源代码, 在 Python 中先获得 URL 对应页面的源代码, 保存到 `htmlSource` 变量中, 然后针对匹配各类 tag 的正则表达式, 分别调用 `re.findall()`, 获得各类 tag 的列表 (因为这个页面中包含的 tag 太多, 每类 tag 只显示前 3 个)。

因为这段程序的输出很多, 在交互式界面下不方便操作和观察, 建议将这些代码单独保存为一个 .py 文件, 比如 `findtags.py`, 然后输入 `python findtags.py` 运行。如果输入 `python` 没有结果 (一般在 Windows 下会出现这种情况), 需要准确设定 PATH 变量, 比如 `d:\Python\python`。之后, 就会看到例 2-9 显示的结果。

#### 例 2-9 使用 `re.findall()` 提取 tag

```
#导入需要的 package
import urllib
import re
#读入 HTML 源代码
sock = urllib.urlopen("http://yahoo.org/")
htmlSource = sock.read()
sock.close()
#匹配, 输出结果 ([0:3]表示取前 3 个)
print "open tags:"
print re.findall(r"<[^>][^>]*[^>]>", htmlSource)[0:3]
print "close tags:"
print re.findall(r"</[>]+>", htmlSource) [0:3]
print "self-closing tags:"
print re.findall(r"<[^>]/+>", htmlSource) [0:3]

open tags:
['<!DOCTYPE html>', '<html lang="en-US" class="y-fp-bg y-fp-pg-grad
bkt701">', '<!-- m2 template 0 -->']
close tags:
['</title>', '</script>', '</script>']
self-closing tags:
['<br/>', '<br/>', '<br/>']
```

## 2.4 点号

上一章讲到了各种字符组，与它相关的还有一个特殊的元字符：点号`.`。一般文档都说，点号可以匹配“任意字符”，点号确实可以匹配“任意字符”，常见的数字、字母、各种符号都可以匹配，如例 2-10 所示。

例 2-10 点号的匹配

```
re.search(r"^.$", "a") != None      # => True
re.search(r"^.$", "0") != None      # => True
re.search(r"^.$", "*") != None      # => True
```

有一个字符不能由点号匹配，就是换行符`\n`。这个字符平时看不见，却存在，而且在处理时并不能忽略（下一章会给出具体的例子）。

如果非要匹配“任意字符”，有两种办法：可以指定使用单行匹配模式，在这种模式下，点号可以匹配合换行符（☞84）；或者使用上一章的介绍“自制”通配字符组`[\s\S]`（也可以使用`[\d\D]`或`[\w\W]`），正好涵盖了所有字符。例 2-11 清楚地说明，这两个办法都可以匹配合换行符。

例 2-11 换行符的匹配

```
re.search(r"^.$", "\n") != None      # => False
#单行模式
re.search(r"(?s)^.$", "\n") != None  # => True
#自制“通配字符组”
re.search(r"^[s\S]$", "\n") != None  # => True
```

## 2.5 滥用点号的问题

因为点号能匹配几乎所有的字符，所以实际应用中许多人图省事，随意使用`.*`或`.*+`，结果却事与愿违，下面以双引号字符串为例来说明。

之前我们使用表达式`"[^"]*"`匹配双引号字符串，而“图省事”的做法是`".*"`。通常这么用是没有问题的，但也可能有意外，例 2-12 就说明了一种如此。

例 2-12 “图省事”的意外结果

```
#字符串的值是"quoted string"
print re.search(r"\.*\"", "\"quoted string\"").group(0)
"quoted string"

#字符串的值是"quoted string" and another"
print re.search(r"\.*\"", "\"quoted string\" and another\"").group(0)
"quoted string" and another"
```

用 `" .* "` 匹配双引号字符串，不但可以匹配正常的双引号字符串 `"quoted string"`，还可以匹配格式错误的字符串 `"quoted string" and another"`。这是为什么呢？

这个问题比较复杂，现在只简要介绍，以说明图省事导致错误的原因，更深入的原因涉及正则表达式的匹配原理，在第8章详细介绍。

在正则表达式 `" .* "` 中，点号 `.` 可以匹配任何字符，`*` 表示可以匹配的字符串长度没有限制，所以 `.*` 在匹配过程结束以前，每遇到一个字符（除去无法匹配的 `\n`），`.*` 都可以匹配，但是到底是匹配这个字符，还是忽略它，将其交给之后的 `"` 来匹配呢？

答案是，具体选择取决于所使用的量词。在正则表达式中的量词分为几类，之前介绍的量词都可以归到一类，叫做**匹配优先量词**（greedy quantifier，也有人翻译为**贪婪量词**<sup>1</sup>）。匹配优先量词，顾名思义，就是在拿不准是否要匹配的时候，优先尝试匹配，并且记下这个状态，以备将来“反悔”。

来看表达式 `" .* "` 对字符串 `"quoted string"` 的匹配过程。

一开始，`"` 匹配，然后轮到字符 `q`，`.*` 可以匹配它，也可以不匹配，因为使用了匹配优先量词，所以 `.*` 先匹配 `q`，并且记录下这个状态【`q` 也可能是 `.*` 不应该匹配的】：

接下来是字符 `u`，`.*` 可以匹配它，也可以不匹配，因为使用了匹配优先量词，所以 `.*` 先匹配 `u`，并且记录下这个状态【`u` 也可能是 `.*` 不应该匹配的】：

.....

现在轮到字符 `g`，`.*` 可以匹配它，也可以不匹配，因为使用了匹配优先量词，所以 `.*` 先匹配 `g`，并且记录下这个状态【`g` 也可能是 `.*` 不应该匹配的】：

最后是末尾的 `"`，`.*` 可以匹配它，也可以不匹配，因为使用了匹配优先量词，所以 `.*` 先匹配 `"`，并且记录下这个状态【`"` 也可能是 `.*` 不应该匹配的】。

这时候，字符串之后已经没有字符了，但正则表达式中还有 `"` 没有匹配，所以只能查询之前保存备用的状态，看看能不能退回几步，照顾 `"` 的匹配。查询到最近保存的状态是：【`"` 也可能是 `.*` 不应该匹配的】。于是让 `.*` “反悔”对 `"` 的匹配，把 `"` 交给 `"`，测试发现正好能匹配，所以整个匹配宣告成功。这个“反悔”的过程，专业术语叫做**回溯**（backtracking），具体的过程如图 2-1 所示。

<sup>1</sup> 许多文档都翻译为“贪婪量词”，单独来看这是没问题的，但考虑到正则表达式中还有其他类型的量词，其英文名字的形式较为统一，所以我在翻译《精通正则表达式》时采用了“匹配优先/忽略优先/占有优先”的名字，也未见读者反对，故此处延用此译法。

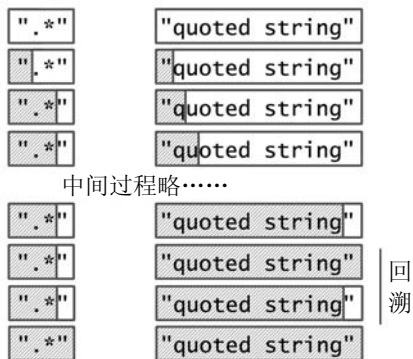


图 2-1 表达式".\*"对字符串"quoted string"的匹配过程



图 2-2 表达式".\*"的匹配过程

如果要准确匹配双引号字符串,就不能图省事使用`.*`,而要使用`"[^"]*"`,过程如图 2-3 所示。

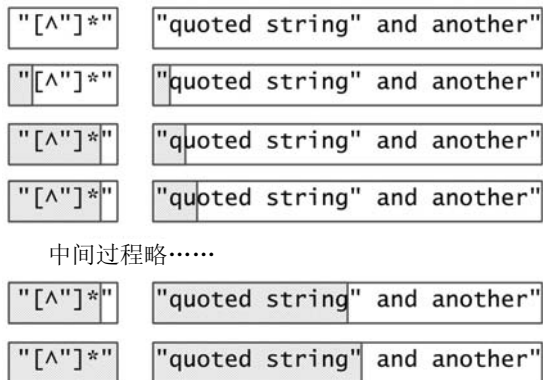


图 2-3 表达式"[^"]\*"的匹配过程

如果把字符串换成`"quoted string" and another`,`.*`会首先匹配第一个双引号之后的所有字符,再进行回溯,表达式中的`"`匹配了字符串结尾的字符`"`,整个匹配宣告完成,过程如图 2-2 所示。

## 2.6 忽略优先量词

也有些时候，确实需要用到 `.*`（或者 `[\s\S]*`），比如匹配 HTML 代码中的 JavaScript 示例就是如此。

```
<script type="text/javascript">...</script>
```

匹配的模式仍然是：匹配 open tag 和 close tag，以及它们之间的内容。open tag 是 `<script type="text/javascript">`，close tag 是 `</script>`，这两段的内容是固定的，非常容易写出对应的表达式，但之间的内容怎么匹配呢？在 JavaScript 代码中，各种字符都可能出现，所以不能用排除型字符组，只能用 `.*`。比如，用一个正则表达式匹配下面这段 HTML 源代码：

```
<script type="text/javascript">
alert("some punctuation <"/");
</script>
```

开头和结尾的 tag 都容易匹配，中间的代码要比较麻烦，因为点号 `.` 不能匹配换行符，所以必须使用 `[\s\S]`（或者 `[\d\D]`、`[\w\W]`）。

```
<script type="text/javascript">[\s\S]*</script>
```

这个表达式确实可以匹配上面的 JavaScript 代码。但是如果遇到更复杂的情况就会出错，比如针对下面这段 HTML 代码，程序运行结果如例 2-13。

```
<script type="text/javascript">
alert("1");
</script>
<br />
<script type="text/javascript">
alert("2");
</script>
```

### 例 2-13 匹配 JavaScript 代码的错误

```
#假设上面的 JavaScript 代码保存在变量 htmlSource 中
jsRegex = r"<script type=\"text/javascript\">[\s\S]*</script>"
print re.search(jsRegex, htmlSource).group(0)
```

```
<script type="text/javascript">
alert("1");
</script>
<br />
```

```
<script type="text/javascript">
alert("2");
</script>
```

用 `<script type="text/javascript">[\s\S]*</script>` 来匹配，会一次性匹配两段 JavaScript 代码，甚至包含之间的非 JavaScript 代码。

按照匹配原理，`[\s\S]*` 先匹配所有的文本，回溯时交还最后的 `</script>`，整个表达式的匹配就成功了，逻辑就是如此，无可改进。而且，这个问题也不能模仿之前双引号字符串匹配，用 `[^"]*` 匹配 `<script...>` 和 `</script>` 之间的代码，因为排除型字符组只能排除单个字符，`^[</script>]` 不能表示“不是 `</script>` 的字符串”。

换个角度来看，通过改变 `[\s\S]*` 的匹配策略解决问题：在不确定是否要匹配的场所，先尝试不匹配的选择，测试正则表达式中后面的元素，如果失败，再退回来尝试 `*` 匹配，如此就没问题了。

循着这个思路，正则表达式中还提供了**忽略优先量词**（lazy quantifier 或 reluctant quantifier，也有人翻译为**懒惰量词**），如果不确定是否要匹配，忽略优先量词会选择“不匹配”的状态，再尝试表达式中之后的元素，如果尝试失败，再回溯，选择之前保存的“匹配”的状态。

对 `[\s\S]*` 来说，把 `*` 改为 `*?` 就是使用了忽略优先量词，`*?` 限定的元素出现次数范围与 `*` 完全一样，都表示“可能出现，也可能不出现，出现次数没有上限”。区别在于，在实际匹配过程中，遇到 `[\s\S]` 能匹配的字符，先尝试“忽略”，如果后面的元素（具体到这个表达式中，是 `</script>`）不能匹配，再尝试“匹配”，这样就保证了结果的正确性，代码见例 2-14。

#### 例 2-14 准确匹配 JavaScript 代码

```
#仍然假设 JavaScript 代码保存在变量 htmlSource 中
jsRegex = r"<script type=\"text/javascript\">[\s\S]*?</script>"
print re.search(jsRegex, htmlSource).group(0)
```

```
<script type="text/javascript">
alert("1");
</script>
```

```
#甚至也可以逐次提取出两段 JavaScript 代码
jsRegex = r"<script type=\"text/javascript\">[\s\S]*?</script>"
for jsCode in re.findall(jsRegex, htmlSource):
    print jsCode + "\n"
```

```
<script type="text/javascript">
alert("1");
</script>
```

```
<script type="text/javascript">
```



```
alert("2");
</script>
```

从表 2-4 可以看到，匹配优先量词与忽略优先量词逐一对应，只是在对应的匹配优先量词之后添加`?`，两者限定的元素能出现的次数也一样，遇到不能匹配的情况同样需要回溯；唯一的区别在于，忽略优先量词会优先选择“忽略”，而匹配优先量词会优先选择“匹配”。

表 2-4 匹配优先量词与忽略优先量词

匹配优先量词	忽略优先量词	限定次数
<code>*</code>	<code>*?</code>	可能不出现，也可能出现，出现次数没有上限
<code>+</code>	<code>+?</code>	至少出现 1 次，出现次数没有上限
<code>?</code>	<code>??</code>	至多出现 1 次，也可能不出现
<code>{m,n}</code>	<code>{m,n}?</code>	出现次数最少为 <i>m</i> 次，最多为 <i>n</i> 次
<code>{m,}</code>	<code>{m,}?</code>	出现次数最少为 <i>m</i> 次，没有上限
<code>{,n}</code>	<code>{,n}?</code>	可能不出现，也可能出现，最多出现 <i>n</i> 次

忽略优先量词还可以完成许多其他功能，典型的例子就是提取代码中的 C 语言注释。

C 语言的注释有两种：一种是在行末，以`//`开头；另一种可以跨多行，以`/*`开头，以`*/`结束。第一种注释很好匹配，使用`//.*`即可，因为点号`.`不能匹配换行符，所以`//.*`匹配的就是从`//`直到行末的文本，注意这里使用了量词`*`，因为`//`可能就是该行最后两个字符；第二种注释稍微复杂一点，因为`/*...*/`的注释和 JavaScript 一样，可能分成许多段，所以必须用到忽略优先量词；同时因为注释可能横跨多行，所以必须使用`[\s\S]`。因此，整个表达式就是`/\*[\s\S]*?\/`（别忘了`*`的转义）。

另一个典型的例子是提取出 HTML 代码中的超链接。常见的超链接形似`<a href="http://somehost/somepath">text</a>`。它以`<a`开头，以`</a>`结束，`href`属性是超链接的地址。我们无法预先判断`<a>`和`</a>`之间到底会出现哪些字符，不会出现哪些字符，只知道其中的内容一直到`</a>`结束<sup>1</sup>，程序代码见例 2-15。

例 2-15 提取网页中所有的超链接 tag

```
#仍然获得 yahoo 网站的源代码，存放在 htmlSource 中
for hyperlink in re.findall(r"<a\s[\s\S]+?</a>", htmlSource):
    print hyperlink
#更多结果未列出
<a href="http://search.yahoo.com/">Web</a>
```

<sup>1</sup> 根据 HTML 规范，`<a>`这个 tag 可用来表示超链接，也可以用作书签，或兼作两种用途，考虑到书签的情况很少见，这里没有做特殊处理。

```
<a href="http://images.search.yahoo.com/images">Images</a>
<a href="http://video.search.yahoo.com/video">Video</a>
```

值得注意的是，在这个表达式中的 `<a` 之后并没有使用普通空格，而是使用字符组简记法 `\s`。HTML 语法并没有规定此处的空白只能使用空格字符，也没有规定必须使用一个空白字符，所以我们用 `\s` 保证“至少出现一个空白字符”（但是不能没有这个空白字符，否则就不能保证匹配 tag name 是 `a`）。

之前匹配 JavaScript 的表达式是 `<script language="text/javascript">[\s\S]*?</script>`，它能应对的情况实在太少了：在 `<script` 之后可能不是空格，而是空白字符；再之后可能是 `type="text/javascript"`，也可能是 `type="application/javascript"`，也可能用 `language` 取代 `type`（实际上 `language` 是以前的写法，现在大都用 `type`），甚至可能没有属性，直接是 `<script>`<sup>1</sup>。

所以必须改造这个表达式，将条件放宽：在 `script` 之后，可能出现空白字符，也可能直接是 `>`，这部分可以用一个字符组 `[\s>]` 来匹配，之后的内容统一用 `[\s\S]+?` 匹配，忽略优先量词保证了匹配进行到最近的 `</script>` 为止。最终得到的表达式就是 `<script[\s>][\s\S]+?</script>`。

对这个表达式稍加改造，就可以写出匹配类似 tag 的表达式。在解析页面时，常见的需求是提取表格中各行、各单元（cell）的内容。表格的 tag 是 `<table>`，行的 tag 是 `<tr>`，单元的 tag 是 `<td>`，所以，它们可以分别用下面的表达式匹配，请注意其中的 `[\s>]`，它兼顾了可能存在的其他属性（比如 `<table border="1">`），同时排除了可能的错误（比如 `<table>`）。

匹配 table	<code>&lt;table[\s&gt;][\s\S]+?&lt;/table&gt;</code>
匹配 tr	<code>&lt;tr[\s&gt;][\s\S]+?&lt;/tr&gt;</code>
匹配 td	<code>&lt;td[\s&gt;][\s\S]+?&lt;/td&gt;</code>

在实际的 HTML 代码中，`table`、`tr`、`td` 这三个元素经常是嵌套的，它们之间存在着包含关系。但是，仅仅使用正则表达式匹配，并不能得到“某个 `table` 包含哪些 `tr`”、“某个 `td` 属于哪个 `tr`”这种信息。此时需要像例 2-16 的那样，用程序整理出来。

### 例 2-16 用正则表达式解析表格

```
# 这里用到了 Python 中的三重引号字符串，以便字符串跨越多行，细节可参考第 14 章
htmlSource = """<table>
<tr><td>1-1</td></tr>
<tr><td>2-1</td><td>2-2</td></tr>
```

<sup>1</sup> 严格说起来，如果只出现 `<script>`，无法保证这里出现的就是 JavaScript 代码，也可能是 VBScript 代码，但考虑到真实世界中的情况，基本可以认为 `<script` 标识的“就是”JavaScript 代码，所以这里不作区分。

```

</table>""

for table in re.findall(r"<table[\s>][\s\S]+?</table>", htmlSource):
    for tr in re.findall(r"<tr[\s>][\s\S]+?</tr>", table):
        for td in re.findall(r"<td[\s>][\s\S]+?</td>", tr):
            print td,
            #输出一个换行符，以便显示不同的行
            print ""
<td>1-1</td>

<td>2-1</td> <td>2-2</td>

```

注：因为 tag 是不区分大小写的，所以如果还希望匹配大写的情况，则必须使用字符组，`table` 写成 `[tT][aA][bB][lL][eE]`，`tr` 写成 `[tT][rR]`，`td` 写成 `[tT][dD]`。

这个例子说明，正则表达式只能进行纯粹的文本处理，单纯依靠它不能整理出层次结构；如果希望解析文本的同时构建层次结构信息，则必须将正则表达式配合程序代码一起使用。

回过头想想双引号字符串的匹配，之前使用的正则表达式是 `"[^"]*"`，其实也可以使用忽略优先量词解决 `".*?"`（如果双引号字符串中包含换行符，则使用 `"[\s\S]*?"`）。两种办法相比，哪个更好呢？

一般来说，`"[^"]*"` 更好。首先，`[^"]` 本身能够匹配合换行符，涵盖了点号 `.` 可能无法应付的情况，出于习惯，很多人更愿意使用点号 `.` 而不是 `[\s\S]`；其次，匹配优先量词只需要考虑自己限定的元素能否匹配即可，而忽略优先量词必须兼顾它所限定的元素与之后的元素，效率自然大大降低，如果字符串很长，两者的速度可能有明显的差异。

而且，有些情况下确实必须用到匹配优先量词，比如文件名的解析就是如此。UNIX/Linux 下的文件名类似这样 `/usr/local/bin/python`，它包含两个部分：路径是 `/usr/local/bin/`；真正的文件名是 `python`。为了在 `/usr/local/bin/python` 中解析出两个部分，使用匹配优先量词是非常方便的。从字符串的起始位置开始，用 `.*` 匹配路径，根据之前介绍的知识，它会回溯到最后（最右）的斜线字符 `/`，也就是文件名之前；在字符串的结尾部分，`[^/]*` 能匹配的就是真正的文件名。前一章介绍过 `^` 和 `$`，它们分别表示“定位到字符串的开头”和“定位到字符串的结尾”，所以应该把 `^` 加在匹配路径的表达式之前，得到 `^.*`，而把 `$` 加在匹配真正文件名的表达式之后，得到 `[^/]*$`，代码见例 2-17。

#### 例 2-17 用正则表达式拆解 Linux/UNIX 的路径

```

print re.search(r"^.*/", "/usr/local/bin/python").group(0)
/usr/local/bin

print re.search(r"[^/]*$", "/usr/local/bin/python").group(0)
python

```

Windows 下的路径分隔符是\，比如 C:\Program Files\Python 2.7.1\python.exe，所以在正则表达式中，应该把斜线字符/换成反斜线字符\。因为在正则表达式中反斜线字符\是用来转义其他字符的，为了表示反斜线字符本身，必须连写两个反斜线，所以两个表达式分别改为`^.*\\`和`[^\\]*$`，代码见例 2-18。

例 2-18 用正则表达式拆解 Windows 的路径

```
#反斜线\必须转义写成\\
print re.search(r"^.*\\", "C:\\Program Files\\Python 2.7.1\\python.exe").group(0)
C:\Program Files\Python 2.7.1\

print re.search(r"[^\\]*$", "C:\\Program Files\\Python 2.7.1\\python.exe").group(0)
python.exe
```

2.7 转义

前面讲解了匹配优先量词和忽略优先量词，现在介绍量词的转义<sup>1</sup>。

在正则表达式中，`*`、`+`、`?`等作为量词的字符具有特殊意义，但有些情况下只希望表示这些字符本身，此时就必须使用转义，也就是在它们之前添加反斜线\。

对常用量词所使用的字符`+`、`*`、`?`来说，如果希望表示这三个字符本身，直接添加反斜线，变为`\+`、`\*`、`\?`即可。但是在一般形式的量词`{m,n}`中，虽然具有特殊含义的字符不止一个，转义时却只需要给第一个`{`添加反斜线即可，也就是说，如果希望匹配字符串`{m,n}`，正则表达式必须写成`\{m,n\}`。

另外值得一提的是忽略优先量词的转义，虽然忽略优先量词也包含不只一个字符，但是在转义时却不像一般形式的量词那样，只转义第一个字符即可，而需要将两个量词全部转义。举例来说，如果要匹配字符串`*?`，正则表达式就必须写作`\*\?`，而不是`\*?`，因为后者的意思是“`*`这个字符可能出现，也可能不出现”。

表 2-5 列出了常用量词的转义形式。

表 2-5 各种量词的转义

量词	转义形式
{n}	\{n}
{m,n}	\{m,n}
{m,}	\{m,}

<sup>1</sup> Java 等语言还支持“占有优先量词 (possessive quantifier)”，但这种量词较复杂，使用也不多，所以本书中不介绍占有优先量词。

( 续表 )

量词	转义形式
{,n}	\{,n}
*	\*
+	\+
?	\?
*?	\*\?
+?	\+\?
??	\?\?

之前还介绍了点号`.`，所以还必须讲解点号的转义：点号`.`是一个元字符，它可以匹配除换行符之外的任何字符，所以如果只想匹配点号本身，必须将它转义为`\.`。

因为未转义的点号可以匹配任何字符，其中也可以包含点号，所以经常有人忽略了对点号的转义。如果真的这样做了，在确实需要严格匹配点号时就可能出错，比如匹配小数（如 `3.14`）、IP 地址（如 `192.168.1.1`）、E-mail 地址（如 `someone@somehost.com`）。所以，如果要匹配的文本包含点号，一定不要忘记转义正则表达式中的点号，否则就有可能出现例 2-19 那样的错误。

例 2-19 忽略转义点号可能导致错误

```
#错误判断浮点数
print re.search(r"^\d+.\d+$", "3.14") != None      # => True
print re.search(r"^\d+.\d+$", "3a14") != None      # => True
#准确判断浮点数
print re.search(r"^\d+\.\d+$", "3.14") != None     # => True
print re.search(r"^\d+\.\d+$", "3a14") != None     # => False
```

# 第 3 章 括号

## 3.1 分组

用字符组和量词可以匹配引号字符串，也可以匹配 HTML tag，如果需要用正则表达式匹配身份证号码，依靠字符组和量词能不能做到呢？

身份证号码是一个长度为 15 或 18 个字符的字符串，如果是 15 位，则全部由数字组成，首位不能为 0；如果是 18 位，则前 17 位全部是数字，末位可能是数字，也可能是 **x**<sup>1</sup>。规则非常明确，可以尝试编写正则表达式了。

首位是数字，不能为 0	[1-9]
除去首末 2 位，剩下 13 位或 16 位，且都是数字	\d{13, 16}
末位可能是数字，也可能是 x	[0-9x]

整个表达式是 `[1-9]\d{13,16}[0-9x]`，它的匹配如例 3-1 所示。

例 3-1 身份证号码的匹配

```
idCardRegex = r"[1-9]\d{13,16}[0-9x]$"
re.search(idCardRegex, "110101198001017032") != None # => True
re.search(idCardRegex, "1101018001017016") != None # => True
re.search(idCardRegex, "11010119800101701x") != None # => True
```

看来，果然能够匹配各种形式的身份证号码，应该没问题。不过这还不够，这个正则表达式应该保证身份证号码的字符串能够匹配，其他字符串不能够匹配，例 3-2 展示了非身份证号码的匹配情况。

例 3-2 身份证号码的错误匹配

```
re.search(idCardRegex, "1101011980010176") != None # => True
re.search(idCardRegex, "110101800101701x") != None # => True
```

<sup>1</sup> 一般来说，最后的 x 可以是小写也可以是大写，但也有些部门规定身份证号码最后的 x 必须是大写 X，这里为讲解方便，只考虑了小写 x 的情况；如果要兼容大写 X 或保证只能出现大写 X，只需要修改最后的字符组 `[0-9x]` 即可。

这两个字符串分明不是身份证号码（第一个有 16 位长，第二个虽然有 15 位长，但末尾是 x），却都匹配了。这是为什么呢？仔细观察所用的正则表达式，会发现两点原因：第一，`\d{13,16}`表示除去首尾两位，中间的部分长度可能在 13~16 之间，而不是“长度要么为 13，要么为 16”；第二，最后的`[0-9x]`只应该对应 18 位身份证号码的情况，但是在这个表达式中，它也可以对应到 15 位身份证号码，而 15 位身份证号码的末位是不能为 x 的！

虽然字符串的长度是可变的，但是除去第一位和最后一位，中间部分的长度必须明确指定，只能是 13 或者 16，而不能使用量词`{13,16}`；另外，末尾一位到底是`[0-9]`（也就是`\d`）还是`[0-9x]`，取决于长度——如果长度是 15 位，则是`\d`；如果长度是 18 位，则是`[0-9x]`。区分两种情况分别考虑，要更加清楚一些。

15 位身份证号码	<code>[1-9]\d{14}</code>
18 位身份证号码	<code>[1-9]\d{14}\d{2}[0-9x]</code>

看来，只要以 15 位号码的匹配为基础，末尾加上可能出现的`\d{2}[0-9x]`即可。这里的`\d{2}[0-9x]`必须作为一个整体，或许不出现（15 位号码），或许出现（18 位号码）。量词`?`可以表示“不出现，或者出现 1 次”，正好用在这里。

但是，在正则表达式`\d{2}[0-9x]?`中，量词`?`只能限定`[0-9x]`的出现，而`\d{2}?[0-9x]?`则更奇怪——即使只出现一个`[0-9x]`，也可以匹配。到底怎样才能把`\d{2}[0-9x]`作为一个整体呢？

答案是：使用括号`(...)`，把正则表达式改写为`[1-9]\d{14}(\d{2}[0-9x])?`。上一章提到过，量词限定之前元素的出现，这个元素可能是一个字符，也可能是一个字符组，还可能是一个表达式——如果把一个表达式用括号包围起来，这个元素就是括号里的表达式，括号内的表达式通常被称为“子表达式”。所以，`(\d{2}[0-9x])?`就表示子表达式`\d{2}[0-9x]`作为一个整体，或许不出现，或许最多出现一次。从例 3-3 可以看到，这个表达式确实可以准确匹配身份证号码。

例 3-3 身份证号码的准确匹配

```
idCardRegex = r"^[1-9]\d{14}(\d{2}[0-9x])?$"
#应该匹配的
re.search(idCardRegex, "110101198001017016") != None # => True
re.search(idCardRegex, "1101018001017016") != None   # => True
re.search(idCardRegex, "11010119800101701x") != None # => True
#不应该匹配的
re.search(idCardRegex, "1101011980010176") != None   # => False
re.search(idCardRegex, "110101800101701x") != None   # => False
```

注：为了方便讲解，我们在正则表达式的两端添加了`^`和`$`，它们分别定位到字符串的起始位置和结束

位置，这样确保了表达式不会只匹配字符串的某个子串；如果要用表达式来提取数据，应当去掉`^`和`$`。下面的例子都遵循这条规则。

括号的这种功能，叫做**分组**（grouping）。如果用量词限定出现次数的元素不是字符或者字符组，而是几个字符甚至表达式，就应该用括号将它们“分为一组”。比如，希望字符串 `ab` 重复出现一次以上，就应该写作 `(ab)+`，此时 `(ab)` 成为一个整体，由量词`+`来限定；如果不用括号而直接写 `ab+`，受`+`限定的就只有 `b`。例 3-4 显示了有括号与无括号的表达式的匹配异同。

#### 例 3-4 用括号改变量词的作用元素

```
re.search(r"^ab+$", "ab") != None      # => True
re.search(r"^ab+$", "abb") != None     # => True
re.search(r"^ab+$", "abab") != None    # => True

re.search(r"^(ab)+$", "ab") != None    # => True
re.search(r"^(ab)+$", "abb") != None   # => False
re.search(r"^(ab)+$", "abab") != None  # => True
```

有了分组，就可以准确表示“长度只能是 *m* 或 *n*”。比如在上面匹配身份证号码的例子中，要匹配一个长度为 13 或者 16 的数字字符串。常犯的错误是使用表达式 `\d{13,16}`，看起来没问题，但长度为 14 或 15 的数字字符串同样会匹配。真正准确的做法是：首先匹配长度为 13 的数字字符串，然后匹配可能出现的长度为 3 的数字字符串，正则表达式就成了 `\d{13}(\d{3})?`。

分组是非常有用的功能，因为使用正则表达式时经常会遇到并没有直接相连，但确实存在联系的部分，分组可以把这些概念上相关的部分“归拢”到一起，以免割裂，下面来看几个例子。

上一章使用表达式 `<[^\s/][^>]*>` 匹配 HTML 中的 open tag，比如 `<table>`，但是这个表达式会匹配 self-closing tag，比如 `<br />`。如果把表达式改为 `<[^\s/][^>]*[^\s/]>`，确实可以避免匹配 self-closing tag，但是因为两个排除型字符组要匹配两个字符，这个表达式又会放过 `<u>` 之类的 open tag，仅仅依靠字符组和量词无法配合解决问题，必须用到括号的分组功能。

`<[^\s/][^>]*[^\s/]>` 错过的只有一种情况，就是 tag name 为单个字母的情况。如果 tag name 不是单个字母，则第一个字母之后，必然会出现这样一个字符串：其中不包含`>`，结尾的字符并不是`/`。最后，才是 tag 结尾的`>`。像图 3-1 所示那样，将这几个元素拆开，能看得更清楚点。

```
tag:    <  u                                >
tag:    <  t      able                      >

表达式:  <  [^\s/]  [^/]*[^\s/]  >
          必须出现  可选出现  必须出现
```

图 3-1 open tag 的准确匹配

所以，需要用一個括号将可选出现的部分分组，再用量词`?`限定，就可以得到兼顾这两种情



况，准确匹配 open tag 的正则表达式了，程序代码如例 3-5 所示。

### 例 3-5 准确匹配 open tag

```
openTagRegex = r"^(^/)([>]*^/)?>$"
re.search(openTagRegex, "<u>") != None      # => True
re.search(openTagRegex, "<table>") != None    # => True

re.search(openTagRegex, "<u/>") != None      # => False
re.search(openTagRegex, "</table>") != None  # => False
```

再看个更复杂的例子。在 Web 服务中，经常并不希望暴露真正的程序细节，所以用某种模式的 URL 来掩盖。比如这个 URL：`/foo/bar_qux.php`，看起来是访问一个 PHP 页面，其实完全不是这样。真正的结构如图 3-2 所示，`foo` 是模块的名称，`bar` 是控制器的名字，`qux` 则是方法名，三个名字都只能出现小写字母。

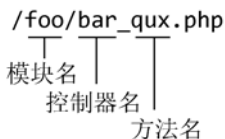


图 3-2 URL 的结构

希望能处理的情况有三种，其他情况都不予考虑。

只有模块名	/foo
只有模块名和控制器名	/foo/bar.php
模块名、控制器名和方法名同时出现	/foo/bar_qux.php

为编写通用的正则表达式来匹配，许多人是这么总结的。

最开始的斜线是必须出现的	/
模块名是必须出现的	[a-z]+
模块名之后的斜线是不一定出现的	/?
控制器名是不一定出现的	[a-z]*
控制器名之后的下划线是不一定出现的	_?
方法名也是不一定出现的	[a-z]*
最后的.php 是不一定出现的	(\.php)?

所以正则表达式就是: `/[a-z]+/?[a-z]*_?[a-z]*(\.php)?。`

仔细看看这个表达式，无论是`/foo`，还是`/foo/bar.php`，抑或是`/foo/bar_qux.php`，都可以匹配，看起来确实没有问题。

可是，这个表达式中只有 `/[a-z]+` 是必须出现的，其他部分都是“不一定出现”——也就是说，其中任意一个或几个部分出现，这个表达式都可以匹配。所以，`/foo/` 也是可以匹配的，`/foo.php` 也是可以匹配的，如例 3-6 所示。

例 3-6 URL 匹配的表达式这里有个下画线

```
urlPatternRegex = r"^/[a-z]+/?[a-z]*_[a-z]*(\\.php)?$"
re.search(urlPatternRegex, "/foo") != None           # => True
re.search(urlPatternRegex, "/foo/bar.php") != None   # => True
re.search(urlPatternRegex, "/foo/bar_qux.php") != None # => True
re.search(urlPatternRegex, "/foo/_") != None         # => True
re.search(urlPatternRegex, "/foo.php") != None       # => True
```

之所以会乱套，根源在于有些元素虽然是“不一定出现”的。可是，“不一定出现”的元素之间却是有关联的：“不一定出现”的元素虽然没有直接相连，却是“要么同时出现，要么同时不出现”的关系。这时候就要梳理清楚逻辑关系，用括号的分组功能把各种分支情况归拢到一起。

`/foo` 是必须出现的；之后存在两种可能：`/bar.php` 或者 `/bar_qux.php`。前一种情况中，开头的 `/`、控制器名 `bar`、结尾的 `.php` 是必须出现的；在后一种情况中，开头的 `/`、控制器名 `bar`、下画线 `_`、模块名 `qux`、结尾的 `.php` 是必须出现的。

/bar.php 对应的表达式	<code>/[a-z]+\.</code> <u><code>php</code></u>
/bar_qux.php 对应的表达式	<code>/[a-z]+</code> <u><code>[a-z]+\.</code></u> <code>php</code>

仔细观察这两个表达式，会发现它们可以合并：把第二个表达式中多出的部分，继续用分组括号配合量词 `?` 表示，塞到第一个表达式中，得到的表达式配合量词 `?`，再加上最开头“必须出现”的 `/foo`，最后得到完整的表达式。

/bar.php 对应的表达式	<code>/[a-z]+\.</code> <u><code>php</code></u>
/bar_qux.php 对应的表达式	<code>/[a-z]+</code> <u><code>[a-z]+\.</code></u> <code>php</code>
合并之后的表达式	<code>/[a-z]+(</code> <u><code>_[a-z]+\.</code></u> <code>)?</code> <code>php</code>
再配合量词?	<code>(/[a-z]+(</code> <u><code>_[a-z]+\.</code></u> <code>)?)?</code>
加上/foo 之后的表达式	<code>/[a-z]+(</code> <u><code>(/[a-z]+(</code></u> <code>_[a-z]+\.</code> <u><code>)?</code></u> <code>)?</code>

从例 3-7 可以看到，这个表达式确实杜绝了错误的匹配。

例 3-7 杜绝了错误匹配的表达式

```
urlPatternRegex = r"^/[a-z]+(/[a-z]+(
```

```
re.search(urlPatternRegex, "/foo/_") != None      # => False
re.search(urlPatternRegex, "/foo.php") != None    # => False
```

关于括号的分组功能，最后来看 E-mail 地址的匹配：E-mail 地址以@分隔为两段，之前的是用户名（username），之后的是主机名（hostname），用户名一般只容许出现数字和字母（现在有些邮件服务商也容许用户名中出现点号等字符了，这种情况复杂些，此处不做考虑），而主机名则是类似 `mail.google.com`、`mail.163.com` 之类的字符串。

用户名的匹配非常简单，其中能出现的字符主要有大写字母 `[A-Z]`、小写字母 `[a-z]`、阿拉伯数字字符 `[0-9]`，下画线 `_`、点号 `.`，所以总的字符组就是 `[A-Za-z0-9_.]`，又可以简化为 `[\w.]`；另一方面，用户名的最大长度是 64 个字符，所以匹配用户名的正则表达式就是 `[\w.]{0,64}`。

主机名匹配的情况则要麻烦一些，简单的情况比如 `somehost.com`；复杂的情况则还包括子域名，比如 `mail.somehost.net`，而且子域名可能不只一级，比如 `mail.sub.somehost.net`。查阅规范可知，主机名被点号分隔为若干段，叫做域名字段（label），每个域名字段中能出现的字符是字母字符、数字字符和横线字符，长度必须在 1~63 之间。下面看几个例子，尝试从中找到主机名的规律。

somehost.net
sub.somehost.net
mail.sub.somehost.net

看来规律是这样的：最后的域名字段是顶级域名<sup>1</sup>，之前的部分可以看作某种模式的重复：该模式由域名字段和点号组成，域名字段在前，点号在后。比如 `somehost.com` 就可以这么看：顶级域名是 `com`，之前是 `somehost.`；`sub.somehost.net` 就可以这么看：顶级域名是 `net`，之前是 `sub.`和 `somehost.`。

匹配域名字段的表达式是 `[-a-zA-Z0-9]{1,63}`，匹配点号的表达式是 `\.`，使用括号的分组功能，把这两个表达式分为一组，用量词 `*` 限定表示“不出现，或出现多次”，就得到匹配主机名的表达式 `([-a-zA-Z0-9]{1,63}\.)*[-a-zA-Z0-9]{1,63}`（因为顶级域名也是一个域名字段，所以即便主机名是 `localhost`，也可以由最后那个匹配域名字段的表达式匹配）。

将匹配用户名的表达式、`@`符号、匹配主机名的表达式组合起来，就得到了完整的匹配 E-mail 地址的表达式：`[-\w.]{0,64}@([-a-zA-Z0-9]{1,63}\.)*[-a-zA-Z0-9]{1,63}`，这个表达式的匹配情况如例 3-8 所示。

<sup>1</sup> 虽然常见的顶级域名是 `cn`、`com`、`info` 之类的，最长 4 个字母，最短 2 个字母，但情况并非都是如此，在有些内部系统中，主机名中并不包含点号，可以视为顶级域名，所以这里认为顶级域名也是一个普通域名字段。

例 3-8 完整匹配 E-mail 地址的正则表达式

```
emailRegex = r"^[ -\w.]{0,64}@[ -a-zA-Z0-9]{1,63}\.)*[ -a-zA-Z0-9]{1,63}$"
#应该匹配的
re.search(emailRegex, "abc@somehost") != None      # => True
re.search(emailRegex, "abc@somehost.com") != None  # => True
re.search(emailRegex, "abc@some-host.com") != None # => True
re.search(emailRegex, "123@somehost.info") != None # => True
re.search(emailRegex, "abc123@somehost.info") != None # => True
re.search(emailRegex, "abc123@sub.somehost.com") != None # => True
re.search(emailRegex, "abc123@m-s.sub.somehost.com") != None # => True
#不应该匹配的
re.search(emailRegex, "abc@.somehost.com") != None # => False
re.search(emailRegex, "a#bc@some-host.commnication") != None # => False
```

3.2 多选结构

之前用表达式 `[1-9]\d{14}(\d{2}[0-9x])?` 匹配身份证号,思路是把 18 位号码多出的 3 位“合并”到匹配 15 位号码的表达式中。能不能直接分情况处理呢? 15 位身份证号就是 `[1-9]` 开头,之后是 14 位数字; 18 位身份证号就是 `[1-9]` 开头,之后是 16 位数字,最后是 `[0-9x]`。只要两个表达式中的一个能够匹配,就是合法的身份证号,这样的思路更加清晰。

15 位身份证号码	<code>[1-9]\d{14}</code>
18 位身份证号码	<code>[1-9]\d{14}\d{2}[0-9x]</code>

答案是可以的,而且仍然使用括号解决问题,只是要用到括号的另一个功能:多选结构 (alternative)。

多选结构的形式是 `(...|...)`,在括号内以竖线 `|` 分隔开多个子表达式,这些子表达式也叫**多选分支** (option); 在一个多选结构内,多选分支的数目没有限制。在匹配时,整个多选结构被视为单个元素,只要其中某个子表达式能够匹配,整个多选结构的匹配就成功;如果所有子表达式都不能匹配,则整个多选结构匹配失败。

回到身份证号码匹配的例子,既然可以区分 15 位和 18 位两种情况,就可以将每种情况对应的表达式作为一个分支,使用多选结构 `([1-9]\d{14}|[1-9]\d{14}\d{2}[0-9x])`。这个表达式的匹配如例 3-9 所示,它同样可以准确验证身份证号码。

例 3-9 用多选结构匹配身份证号码

```
idCardRegex = r"^[1-9](\d{14}|\d{14}\d{2}[0-9x])$"
#应该匹配的
re.search(idCardRegex, "110101198001017016") != None # => True
```

```
re.search(idCardRegex, "1101018001017016") != None    # => True
re.search(idCardRegex, "11010119800101701x") != None    # => True
#不应该匹配的
re.search(idCardRegex, "1101011980010176") != None    # => False
re.search(idCardRegex, "110101800101701x") != None    # => False
```

多选结构在实际中经常用到，匹配 IP 地址就是如此：IP 地址（暂不考虑 IPv6）分为四段（四个字节），每段都是八位二进制数，换算成常见的十进制，取值在 0~255 之间，中间以点号.分隔。点号.的匹配非常容易，用\.就可以，所以暂且忽略它，只考虑匹配这个数值的问题，而且因为 4 段 IP 地址的取值范围是相同的，只考虑其中一段的匹配即可。

要匹配十进制形式的 IP 地址，最常见的正则表达式就是 `[0-9]{1,3}`，也就是 1~3 位十进制数字。粗看起来，这个表达式没什么错，细看却有很大问题。因为 256、999 这样的数值，显然不在 0~255 之间，却可以由 `[0-9]{1,3}` 匹配。

细致一点的表达式似乎是 `[0-2][0-5][0-5]`，这样就限制了数值只能是在 255 以内……不过，仔细想想，因为限定了第二位（十位）和第三位（个位）都只能出现 0~5 之间的字符，表达式没法匹配 168 之类的数值。

其实，问题可以这样解决：先用表达式匹配这个字符串，再将它转换为整数类型的变 `x`，判断 `x` 是否在 0 和 255 之间：`0<=x && x<=255`。没错，这确实是一个解决问题的思路，只是麻烦一点，最好能用正则表达式“一次性”搞定这个问题。仔细想想就会发现，正则表达式虽然直接表示“匹配一段数值在 0~255 之间的文本”，但可以分几种情况描述符合这样规则的文本。

如果是 1 位数，那么，对数字没有限制	<code>[0-9]</code>
如果是 2 位数，那么，对数字没有限制	<code>[0-9]{2}</code>
如果是 3 位数	
如果第 1 位数字是 1，那么第 2、3 位数字都没有限制	<code>1[0-9][0-9]</code>
如果第 1 位数字是 2	
如果第 2 位数字是 0~4，那么第 3 位数字没有限制	<code>2[0-4][0-9]</code>
如果第 2 位数字是 5，那么第 3 位数字只能是 0~5	<code>25[0-5]</code>

虽然不如 `0<=x && x<=255` 的判断简便，但如果文本符合其中任何一条规则（或者说，只要其中任何一个正则表达式能匹配），就可以判断它“表示数字的数值在 0~255 之间”。用多选结构把这几条规则对应的表达式合并起来，就得到了表达式 `([0-9]| [0-9]{2}| 1[0-9][0-9]| 2[0-4][0-9]| 25[0-5])`，它的匹配如例 3-10 所示。

例 3-10 准确匹配 0~255 之间的字符串

```
partRegex = r"^([0-9]| [0-9]{2}| 1[0-9][0-9]| 2[0-4][0-9]| 25[0-5])$"
#应该匹配的
```

```
re.search(partRegex, "0") != None # => True
re.search(partRegex, "98") != None # => True
re.search(partRegex, "168") != None # => True
#不应该匹配的
re.search(partRegex, "256") != None # => False
```

如果要更完善一点，能识别 030、005 这样的数值，可以修改对应的子表达式，为一位数和两位数的情况增加之前可能出现 0 的匹配，得到表达式 `((00)?[0-9]|0?[0-9]{2})|1[0-9][0-9]|2[0-4][0-9]|25[0-5]`。

上面讲解的，其实是用正则表达式匹配数值在某个范围内的字符串的通用模式，它很重要，因为许多时候会遇到类似的任务，比如匹配月（1~12）、日（不考虑只有 30 天的情况，粗略记为 1~31）、小时（0~24）、分钟（00~60）的正则表达式，用正则表达式解决这类问题，会用到同样的模式。

月	<code>(0?[1-9] 1[012])</code>
日	<code>(0?[1-9] 1[12]\d 3[01])</code>
小时	<code>(0?[1-9] 01\d 2[0-4])</code>
分钟	<code>(0?[1-9] 0[0-5]\d 60)</code>

这个模式还可以用于匹配手机号码：大陆的手机号码是 11 位的，前面 3 位是号段，到目前为止有 130~139 号段、150~153、155~156、180、182、185~189 号段，用多选分支 `(13[0-9]|15[0-356]|18[025-9])` 可以很准确地匹配号段；之后的 8 位一般没有限制，只要是数字即可，用 `\d{8}` 匹配。另外，手机号码最开头可能有 0 或者+86，它可以用 `(0|\+86)` 匹配，因为整个部分是可能出现的，所以需要加上量词，也就是 `(0|\+86)?` 最后得到的正则表达式就是 `(0|\+86)?(13[0-9]|15[0-356]|18[025-9])\d{8}`。

多选结构还可以解决更复杂的问题，比如上一章的 tag 匹配问题，当时使用的表达式是 `<[>]+>`，一般来说，这个表达式是没有问题的，但也有可能 tag 内部还是会出现 `>` 符号，比如 `<input name=txt value=">">`。这类问题使用字符组解决不了，使用多选结构则可以解决。

仔细分析 tag 中可能出现 `>` 它只可能作为属性（attribute）出现在单引号字符串和双引号字符串中，根据 html 规范，引号字符串中不能出现嵌套转义的引号，所以单引号字符串可以用 `'[^']*'` 来匹配，双引号字符串可以用 `"[^"]*"` 来匹配，相应的，其他内容则可以用 `[^'">]` 来匹配，所以更完善的表达式是 `<('[''^']*'|"["^"]*"|"[^"]">|"[^"]">)+>`。它的匹配情况见例 3-11。

例 3-11 准确的 HTML tag 匹配

```
tagRegex = r"^<('[''^']*'|"["^"]*"|"[^"]">|"[^"]">)+>$"
re.search(tagRegex, "<input name=txt value=\">\">") != None # => True
re.search(tagRegex, "<input name=txt value='>'>") != None # => True
re.search(tagRegex, "<a>") != None # => True
```

请注意其中的量词，因为单引号字符串和双引号字符串都可以是空字符串，比如 `alt=""` 或 `alt=''`，所以匹配其中文本的内容使用 `*`；而 `['>']` 则没有使用量词，因为它存在于多选结构内部，多选结构外部有 `+` 量词限制，保证了它不只匹配一个字符。如果在多选结构内部使用 `['>']*`，虽然看来似乎没错，却可能导致非常奇怪的结果，不过现在不用关心，详细的讲解在第 135 页。

关于多选结构，最后还要补充三点。

第一，多选结构的一般表示法是 `(option1|option2)`（其中 `option1` 和 `option2` 是两个作为多选分支的正则表达式），多选结构中一般会同时使用括号 `()` 和竖线 `|`；但是如果没有括号 `()`，只出现竖线 `|`，仍然是多选结构。从例 3-12 可以看到，`ab|cd` 既可以匹配 `ab`，也可以匹配 `cd`。

### 例 3-12 没有括号的多选结构

```
re.search(r"ab|cd", "ab") != None    # => True
re.search(r"ab|cd", "cd") != None    # => True
```

在多选结构中，竖线 `|` 用来分隔多选结构，而括号 `()` 用来规定整个多选结构的范围，如果没有出现括号，则将整个表达式视为一个多选结构，所以 `ab|cd` 等价于 `(ab|cd)`。如果在某些文档中看到没有括号的多选结构，不用奇怪。

不过，我还是推荐明确写出两端的括号，这样更形象，也能避免一些错误。如果你仔细看，就会发现在上面的表达式中，并没有使用 `^` 和 `$` 定位字符串的起始位置和结束位置，按道理说，加上之后应该匹配更加准确，结果却并非如此。

因为竖线 `|` 的优先级很低（关于优先级，☞106），所以 `^ab|cd$` 其实是 `(^ab|cd$)`，而不是 `^(ab|cd)$`，它的真正意思是“字符串开头的 `ab` 或者字符串结尾的 `cd`”，而不是“只包含 `ab` 或 `cd` 的字符串”，代码见例 3-13。

### 例 3-13 没有括号的多选结构

```
re.search(r"^ab|cd$", "abc") != None    # => True
re.search(r"^ab|cd$", "bcd") != None    # => True
re.search(r"^(ab|cd)$", "abc") != None  # => False
re.search(r"^(ab|cd)$", "bcd") != None  # => False
```

第二，多选分支并不等于字符组。多选分支看起来类似字符组，如 `[abc]` 能匹配的字符串和 `(a|b|c)` 一样，`[0-9]` 能匹配的字符串和 `(0|1|2|3|4|5|6|7|8|9)` 一样。从理论上说，可以完全用多选结构来替换字符组，但这种做法并不推荐，理由在于：首先，`[abc]` 比 `(a|b|c)` 要简洁许多，在多选结构中的每个分支都必须明确写出，不能使用 - 范围表示法，

`(0|1|2|3|4|5|6|7|8|9)`比`[0-9]`麻烦很多；其次，大多数情况下，`[abc]`比`(a|b|c)`的效率要高很多。所以，能用字符组解决的问题，最好不用等价的多选结构。

反过来，多选结构不一定能对应到字符组。因为字符组的每个“分支”的长度相同，而且只能是单个字符；而多选结构的每个“分支”的长度没有限制，甚至可以是复杂的表达式，比如`(abc|b+c*ab)`，字符组完全无能为力。

多选分支和字符组的另一点重要区别（同时也是最常犯的错误）是：排除型字符组可以表示“无法由某几个字符匹配的字符”，多选结构却没有对应的结构表示“无法由某几个表达式匹配的字符串”。从例 3-14 可以看到，`[^abc]`表示“匹配除 a、b、c 之外的任意字符”，`(^a|b|c)`却不能表示“匹配除 a、b、c 之外的任意字符串”。

#### 例 3-14 多选结构不能表示“无法由某几个表达式匹配的字符串”

```
re.search(r"^(a|b|c)", "ab") != None      # => True
re.search(r"^(a|b|c)", "cd") != None      # => True
```

在实际开发中确实可能遇到这种需求，不过没有现场的解法。如果你现在就希望匹配“无法由某几个表达式匹配的字符串”，请参考第 140 页。

第三，多选分支的排列是有讲究的。比如这个表达式`(jeff|jeffrey)`，用它匹配 `jeffrey`，结果到底是 `jeff` 还是 `jeffrey` 呢？这个问题并没有标准的答案，本书介绍的 Java、.NET、Python、Ruby、JavaScript、PHP 中，多选结构都会优先选择最左侧的分支。这一点从例 3-15 看得很清楚：如果使用字符串是 `jeffrey`，正则表达式是`(jeff|jefferey)`还是`(Jeffrey|jeff)`，结果是不一样的（此处仅以 Python 为例，本书中介绍的其他语言中的结果与此相同）。

#### 例 3-15 多选结构的匹配顺序

```
print re.search(r"(jeffrey|jeff)", " jeffrey").group(0)
jeffrey

print re.search(r"(jeff|jeffrey)", " jeffrey").group(0)
jeff
```

在实际开发中可能会遇到这样的情况：统计一段文本中，“湖南”和“湖南省”分别出现的次数。如果直接查找“湖南”，可能会将“湖南省”中的“湖南”也找出来，如果使用多选结构`(湖南省|湖南)`，就可以一次性找出所有“湖南”和“湖南省”，再按照字符串的长度分别计数，就可以得到两者出现的次数了。

不过，`(湖南省|湖南)`只是一个针对特殊应用的例子。在平时使用中，如果出现多选结构，应当尽量避免多选分支中存在重复匹配，因为这样会大大增加回溯的计算量。也就是说，应当避免这样的情况：针对多选结构`(option1|regex2)`，某段文本既可以由 `option1` 匹配，也可



以由 `option2` 匹配。如果出现了这样的多选结构，效率可能会受到极大影响（第 160 页总结了可能影响效率的几种写法），尤其在受量词限定的多选结构中更是如此：一般人不会遇到 `(a|[ab])` 这类多选结构，但 `([0-9]|\w)` 之类则一不留神就会遇到。

### 3.3 引用分组

括号不仅仅能把有联系的元素归拢起来并分组，还有其他的作用——使用括号之后，正则表达式会保存每个分组真正匹配的文本，等到匹配完成后，通过 `group(num)` 之类的方法“引用”分组在匹配时捕获的内容（这个方法之前已经出现过）。其中，`num` 表示对应括号的编号，括号分组的编号规则是从左向右计数，从 1 开始。因为“捕获”了文本，所以这种功能叫做**捕获分组**（capturing group）。对应的，这种括号叫做捕获型括号。

举个例子，我们经常遇到诸如 `2010-12-22`、`2011-01-03` 这类表示日期的字符串，希望从中提取出年、月、日之类的信息，就可借助捕获分组来实现。正则表达式中，每个捕获分组都有一个编号，具体情况如图 3-3 所示。

字符串：	<code>2010</code>	-	<code>12</code>	-	<code>22</code>
字符串：	<code>2011</code>	-	<code>01</code>	-	<code>03</code>
表达式：	<code>(\d{4})</code>	-	<code>(\d{2})</code>	-	<code>(\d{2})</code>
分组编号：	1		2		3

图 3-3 分组及编号

一般来说，正则表达式匹配完成之后，都会得到一个表示“匹配结果”的对象，对它调用获取分组的方法，传入分组编号 `num`，就可以得到对应分组匹配的文本。第 1 章介绍过，如果匹配成功，`re.search()` 返回一个 `MatchObject` 对象。如果只需要知道“是否能匹配”，判断它是否为 `None` 即可；但如果获取了 `MatchObject` 对象，也可以通过对应的方法，显示匹配结果的详细信息。使用 `MatchObject.group(num)`，就可以引用正则表达式中编号为 `num` 的分组匹配的文本。从例 3-16 可以看到，通过引用编号为 1、2、3 的捕获分组，分别获得了年、月、日的信息。

#### 例 3-16 引用捕获分组

```
print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(1)
2010

print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(2)
12

print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(3)
22
```

前面说过，*num* 的编号从 1 开始。不过，也有编号为 0 的分组，它是默认存在的，对应整个表达式匹配的文本。在许多语言中，如果调用 `group()` 方法，不给出参数 *num*，默认就等于调用 `group(0)`，比如 Python 就是如此，代码见例 3-17。

#### 例 3-17 默认存在编号为 0 的分组

```
print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(0)
2010-12-22

print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group()
2010-12-22
```

有些正则表达式里可能包含嵌套的括号，比如在上面的例子中，除了能单独提取出年、月、日之外，再给整个表达式加上一重括号，就出现了嵌套括号，这时候括号的编号是怎样的呢？答案很简单：无论括号如何嵌套，分组的编号都是根据开括号出现顺序来计数的；开括号是从左向右数起第多少个开括号，整个括号分组的编号就是多少。图 3-4 举例说明了这种编号规则，具体的代码见例 3-18。

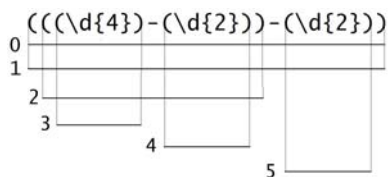


图 3-4 分组编号只取决于开括号出现的顺序

#### 例 3-18 嵌套的括号

```
nestedGroupingRegex = r"((\d{4})-(\d{2}))-(\d{2})"
print re.search(nestedGroupingRegex, "2010-12-22").group(0)
2010-12-22

print re.search(nestedGroupingRegex, "2010-12-22").group(1)
2010-12-22

print re.search(nestedGroupingRegex, "2010-12-22").group(2)
2010-12

print re.search(nestedGroupingRegex, "2010-12-22").group(3)
2010

print re.search(nestedGroupingRegex, "2010-12-22").group(4)
12

print re.search(nestedGroupingRegex, "2010-12-22").group(5)
22
```

上一章用正则表达式 `<a\s[\s\S]+?</a>` 提取 HTML 中的所有的超链接 tag，配合括号的分组功能，可以更进一步，依靠引用分组把超链接的地址和文本分别提取出来。通常的超链接 tag 类似这样：`<a href="url">text</a>`。其中 `url` 是超链接地址，`text` 是文本，为了准确获取这两部分内容，可以把表达式改为 `<a\s+href="([^\"]+)">([^\<]+)</a>`。

其中给匹配 `url` 和 `text` 的表达式分别加上括号，就是 `([^\"]+)` 和 `([^\<]+)`（注意其中 `<a` 之后是 `\s+`，因为这里需要的是空白字符，而不限定是空格字符，而且可能不止一个字符）。

当然这只是最简单的情况，在等号=两端可能还有空白字符，比如 `<a href = "url">text</a>`，所以正则表达式中的 `=` 两端也应该添加 `\s*`，于是得到 `<a\s+href\s*=\s*"([^\"]+)">([^\<]+)</a>`。

不过，属性既可以用双引号字符串表示，也可以用单引号字符串表示，比如 `<a href='url'>text</a>`；甚至可以不用引号，比如 `<a href=url>text</a>`。为了处理这两种情况，需要继续改造表达式：首尾出现的单引号或者双引号字符用 `['"]?` 即可匹配；真正的 URL，既不能包含单引号，也不能包含双引号，还不能是空白字符，所以可以用 `^[^'" \s]+` 匹配，而且这部分是需要提取出来的，别忘了它外面的括号。于是得到了最后的表达式 `<a\s+href\s*=\s*"([^\"]+)"|'([^']*)'|([^\s'"]+)">([^\<]+)</a>`。

现在表达式已经编写完毕，第一个括号内的表达式用来匹配 `url`，第二个括号内的表达式用来匹配 `text`，所以如果要提取 `url` 和 `text`，应该使用编号为 1 和 2 的分组。下面仍然以 `yahoo.com` 的首页为例来看看结果。需要说明的是，如果使用 `re.findall()`，而且正则表达式中出现了捕获型括号，那么返回数组的每个元素都是数组，其中各个元素对应各个分组的文本，所以直接用下标 2 访问得到第二个分组对应的文本，不必显式调用 `group(2)`，代码见例 3-19。

### 例 3-19 用分组提取出超链接的详细信息

```
# yahoo.com 的源代码已经保存在 htmlSource 中
hrefTagRegex = r"<a\s+href\s*=\s*"([^\"]+)"|'([^']*)'|([^\s'"]+)">([^\<]+)</a>"
for hyperlink in re.findall(hrefTagRegex, htmlSource):
    print hyperlink[2], hyperlink[1]

Web      http://search.yahoo.com/
Images   http://images.search.yahoo.com/images
Video    http://video.search.yahoo.com/video
.....更多结果未列出
```

类似的，还可以提取出网页标题（`<head>`）或网页中的图片链接（`<img>`）的表达式<sup>1</sup>。应

<sup>1</sup> 虽然 HTML 4.0 规范中没有限制 tag 名的大小写，但 XHTML 规范要求都使用小写，在平时处理中，tag 名也大多使用小写，如果实在要处理大写的情况，可以使用字符组，比如 `img` 写作 `[iI][mM][gG]`，更省事的办法是指定不区分大小写的模式，☞81。

当注意的是，匹配<img>时，在<img 和 src 之间可能还有其他内容，比如 width=750 之类，所以不能仅仅用\s+匹配，而应当添加[^>]\*?。在 src=...之后也是同样如此。表 3-1 总结了匹配网页标题和图片链接的表达式。

表 3-1 提取网页标题和图片链接的正则表达式

任务	表达式	提取分组编号
网页标题	<head>([^<]+)</head>	1
图片链接	<img\s+[^>]*?src=[\'"]?([^\'\"\\s]+)[\'"]?[^>]*>	1

应当记住的是，引用分组时，引用的是分组对应括号内的表达式捕获的文本。在这个问题上，正则表达式新手常犯错误。例 3-20 仍然是用正则表达式匹配日期字符串，两个表达式能匹配的字符串是完全相同的，引用分组的编号也是相同的，结果却不同。

例 3-20 新手容易弄错分组的结构

```
re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(1)
2010

re.search(r"(\d){4}-(\d{2})-(\d{2})", "2010-12-22").group(1)
0
```

在第一个表达式中，编号为 1 的括号是(\d{4})，其中的\d{4}是“匹配四个数字字符”的子表达式。在第二个表达式中，编号为 1 的括号是(\d)，其中的\d是“匹配一个数字字符”的子表达式，因为之后有量词{4}，所以整个括号作为单个元素，要重复出现 4 次，而且编号都是 1；于是每重复出现一次，就要更新一次匹配结果。所以在匹配过程中，编号为 1 的分组匹配的文本的值，依次是 2、0、1、0，最后的结果是 0。在实际使用时，常常有人忽略了这一细节，得到匪夷所思的匹配结果。

引用分组捕获的文本，不仅仅用于数据提取，也可以用于替换，有时候这么做非常方便。仍然举上面的日期的例子，比如希望将 YYYY-MM-DD 格式的日期变为 MM/DD/YYYY，就可以使用正则表达式替换。

在 Python 语言中进行正则表达式替换的方法是 re.sub(pattern, replacement, string)，其中 pattern 是用来匹配被替换文本的表达式，replacement 是要替换成的文本，string 是要进行替换操作的字符串，比如 re.sub(r"[a-z]", " ", string)就是将 string 中的每一个小写字母替换为一个空格。程序运行结果如例 3-21。

例 3-21 正则表达式替换

```
print re.sub(r"[a-z]", " ", "1a2b3c")
1 2 3
```

在 *replacement* 中也可以引用分组，形式是 `\num`，其中的 *num* 是对应分组的编号。不过，*replacement* 并不是一个正则表达式，而是一个普通字符串。根据字符串中的转义规定，`\t` 表示制表符，`\n` 表示换行符，`\1`、`\2` 却不是字符串中的合法转义序列，所以也必须指定 *replacement* 为原生字符串（☞93）。例 3-22 说明了如何通过 *replacement* 中使用了引用分组，转换日期字符串的格式。

### 例 3-22 在替换中使用分组

```
print re.sub(r"(\d{4})-(\d{2})-(\d{2})", r"\2/\3/\1", "2010-12-22")
12/22/2010

print re.sub(r"(\d{4})-(\d{2})-(\d{2})", r"\1年\2年\3日", "2010-12-22")
2010年12月22日
```

值得注意的是，如果想在 *replacement* 中引用整个表达式匹配的文本，不能使用 `\0`，即使用原生字符串也不行。因为在字符串中，`\0` 开头的转义序列通常表示用八进制形式表示的字符，`\0` 本身表示 ASCII 字符编码为 0 的字符。如果一定要引用整个表达式匹配的文本，则可以稍加变通，给整个表达式加上一对括号，之后用 `\1` 来引用，如例 3-23。

### 例 3-23 在替换中，使用 `\1` 替代 `\0`

```
#ASCII 编码为 0 的字符无法显示
print re.sub("(\\d{4})-(\\d{2})-(\\d{2})", "\\0", "2010-12-22")
print re.sub("(\\d{4})-(\\d{2})-(\\d{2})", r"\0", "2010-12-22")

print re.sub("((\\d{4})-(\\d{2})-(\\d{2}))", "[\\1]", "2010-12-22")
[2010-12-22]

print re.sub("((\\d{4})-(\\d{2})-(\\d{2}))", r"[\1]", "2010-12-22")
2010-12-22
```

## 3.3.1 反向引用

英文的不少单词中都有重叠出现的字母，比如 **shoot** 或 **beep**，如果希望检查某个单词是否包含重叠出现的字母，该怎么办呢？

匹配字母的表达式是 `[a-z]`（这里暂时不考虑大写的情况），所以最先想到的往往是用两个字符组 `[a-z][a-z]` 来匹配，但这样做并不对，因为重叠出现的字母是不确定的。假设字符串是 **at**，**a** 可以由第一个 `[a-z]` 匹配，**t** 可以由第二个 `[a-z]` 匹配，但是因为前一个 `[a-z]` 和后一个 `[a-z]` 之间并没有联系，所以 `[a-z][a-z]` 其实只能匹配两个小写字母，不关心它们是否相同。

这个问题有点复杂。“重叠出现”的字母，取决于第一个 `[a-z]` 在运行时的匹配结果，而不

能预先设定。也就是说必须“知道”之前匹配的确切内容：如果前面的`[a-z]`匹配的是`e`，后面就只能匹配`e`；如果前面的`[a-z]`匹配的是`o`，后面就只能匹配`o`。

前面我们看到了引用分组，能引用某个分组内的子表达式匹配的文本，但引用都是在匹配完成后进行的，能不能在正则表达式中引用呢？

答案是可以的，这种功能被称作**反向引用**（back-reference），它允许在正则表达式内部引用之前的捕获分组匹配的文本（也就是左侧），其形式也是`\num`，其中 *num* 表示所引用分组的编号，编号规则与之前介绍的相同。

根据反向引用，查找连续重叠字母的表达式就是`([a-z])\1`，其中的`[a-z]`匹配第一个字母，再用括号将匹配分组，然后用`\1`来反向引用，这个表达式的匹配情况见例 3-24。

### 例 3-24 用反向引用匹配重复字母

```
re.search(r"^[a-z]\1$", "aa") != None      # => True
re.search(r"^[a-z]\1$", "dd") != None      # => True
re.search(r"^[a-z]\1$", "ac") != None      # => False
```

在日常开发中，我们可能经常需要反向引用来建立前后联系。最常见的例子就是解析 HTML 代码时匹配 tag。之前我们说过，tag 包括 open tag 和 close tag，open tag 和 close tag 经常是成对出现的，比如`<bold>text</bold>`或`<h1>title</h1>`。

有了反向引用功能，就可以先匹配 open tag，再匹配其他内容，直到最近的 close tag 为止：在匹配 open tag 时，用一个括号分组匹配 tag name 的表达式`([<^>]+)`；在匹配 close tag 时，用`\1`引用之前匹配的 tag name，就完成了配对（要注意的是，这里需要用到忽略优先量词`??`，否则可能会出现错误匹配，理由在上一章匹配 JavaScript 代码时讲过）。最后得到的表达式就是`<([<^>]+)>[\s\S]*?</\1>`，这个表达式的匹配如例 3-25 所示。

### 例 3-25 用反向引用匹配成对的 tag

```
pairedTagRegex = r"<([<^>]+)>[\s\S]*?</\1>"
#应该匹配的
re.search(rpairedTagRegex, "<bold>text</bold>") != None      # => True
re.search(rpairedTagRegex, "<h1>title</h1>") != None          # => True
#不应该匹配的
re.search(rpairedTagRegex, "<h1>text</bold>") != None          # => False
```

也有些 tag 更复杂一点，比如`<span class="class1">text</span>`，在 tag 名之后有一个空白字符，然后是其他属性，此时原有的表达式就无法匹配了。为应对这类情况，应当修改表达式让分组 1 准确匹配 tag name，它可以是数字、小写字母、大写字母，所以将它修改为`<([a-zA-Z0-9]+)\s[<^>]+>[\s\S]*?</\1>`，但满足了`\s[<^>]+>`的匹配，就无法应对之前的那些 open tag。为了兼容两种情况，必须用括号分组和量词`?`来限定，改为也就是`(\s[<^>]+)?`，

最后的表达式就是 `<([a-zA-Z0-9]+)(\s[^>]+)?>[\s\S]*?</\1>`。具体程序如例 3-26。

例 3-26 用反向引用匹配更复杂的成对 tag

```
pairedTagRegex = r"<([a-zA-Z0-9]+)(\s[^>]+)?>[\s\S]*?</\1>"
re.search(pairedTagRegex, "<bold>text</bold>") != None      # => True
re.search(pairedTagRegex, "<h1>title</h1>") != None          # => True
re.search(pairedTagRegex, "<span class=\"class1\">text</span>") != None  #=> True
re.search(pairedTagRegex, "<h1>text</bold>") != None          # => False
```

反向引用还可以用在其他很多地方，比如在处理中文文本时，查找“浩浩荡荡”、“清清白白”之类 AABB，或者“如火如荼”、“越快越好”之类 AXAY 的四字词组。

关于反向引用，还有一点需要强调：反向引用重复的是对应捕获分组匹配的文本，而不是之前的表达式；也就是说，反向引用的是由之前表达式决定的具体文本，而不是符合某种规则的未知文本。这一点，新手常犯错误。

仍然以匹配 IP 地址为例，前面说过，IP 地址分 4 段（4 个字节），匹配其中每一段的表达式是 `(0{0,2}[0-9]|0?[0-9]{2}|1[0-9][0-9]|2[0-4][0-9]|25[0-5])`，之间用点号分隔，所以匹配完整 IP 地址的表达式应该是用量词重复这个子表达式，而不是用反向引用重复这个表达式匹配的文本。例 3-27 对比了这两个表达式，其中第二个表达式中使用了反向引用，故而要求后面 3 段与第 1 个字段完全一样，所以它只能匹配 `8.8.8.8` 之类地址，而不能匹配 `192.168.0.1` 之类地址。

例 3-27 匹配 IP 地址的正则表达式

```
#匹配其中一段的表达式
#segment = r"(0{0,2}[0-9]|0?[0-9]{2}|1[0-9][0-9]|2[0-4][0-9]|25[0-5])"
#正确的表达式
ipAddressRegex = r"(" + segment + r"\.){3}" + segment
#错误的表达式
ipAddressRegex = segment + r"\.\1\.\1\.\1"
```

### 3.3.2 各种引用的记法

根据前面的介绍，对分组的引用可能出现在三种场合：在匹配完成后，用 `group(num)` 之类的方法提取数据；在进行正则表达式替换时，用 `\num` 引用；在正则表达式内部，用 `\num` 引用。

不过，这只是 Python 语言的规定，事情并不总是如此：`group(num)` 之类的方法，在各种语言中都是差不多的；但是在有些语言中，替换时引用的记法和正则表达式内部引用的记法是不同的。表 3-2 总结了各种常用语言中的两类记法<sup>1</sup>。

<sup>1</sup> 在很多文档中都表示为 `\n`，但 `\n` 容易误解为换行符，所以本书中统一用 `\num` 表示。

表 3-2 各种语言中引用分组的记法

语言	表达式中的反向引用	替换中的反向引用
.NET	<code>\num</code>	<code>\$num</code>
Java	<code>\num</code>	<code>\$num</code>
JavaScript	<code>\$num</code>	<code>\$num</code>
PHP	<code>\num</code>	<code>\num</code> 或 <code>\$num</code> (PHP 4.0.4 以上版本)
Python	<code>\num</code>	<code>\num</code>
Ruby	<code>\num</code>	<code>\num</code>

看起来`\num`和`$num`差别不大：`\1`或者`$1`表示第 1 个捕获分组，`\2`或者`$2`表示第 2 个捕获分组……不过一般来说，`$num`要好于`\num`。原因在于，`$0`可以准确表示“第 0 个分组（也就是整个表达式匹配的文本）”，而`\0`则不行，因为在不少语言的字符串中，`\num`本身是一个有意义的转义序列，它表示值为 `num` 的 ASCII 字符，所以`\0`会被解释为“ASCII 编码为 0 的字符”。但是反向引用不存在这个问题，因为不能在正则表达式还没匹配结束时，就用`\0`引用整个表达式匹配的文本。

但无论是`\num`还是`$num`，都有可能遇到二义性的问题：如果出现了`\10`（或者`$10`，这里以`\num`为例），它到底表示第 10 个捕获分组`\10`，还是第 1 个捕获分组`\1`之后跟着一个字符 `0`？Python 的结果见例 3-28。

例 3-28 可能具有二义性的反向引用

```
print re.sub(r"(\d)", r"\10", "123")
Traceback (most recent call last):
sre_constants.error: invalid group reference
```

原来`\10`会被解释成“第 10 个捕获分组匹配的文本”，而不是“第 1 个捕获分组匹配的文本之后加上字符 `0`”。如果我们就是希望做到后面这步，Python 提供了`\g<num>`表示法，将`\10`写成`\g<1>0`，这样同时也避免了替换时无法使用`\0`的问题，代码如例 3-29。

例 3-29 使用 `g<n>` 消除二义性

```
print re.sub(r"(\d)", r"\g<1>0", "123")
102030
```

PHP 中也有专门的记法解决这类问题，在替换时可以使用`\${num}`的写法，准确标注所引用分组的编号，也就是说，`\${1}0`表示“第 1 个捕获分组之后加上 `0`”，`${10}`表示“第 10 个捕获分组”。而`$10`，在第 10 个捕获分组存在的情况下，表示该捕获分组；否则，被视为空字符串。PHP 的代码见例 3-30。



## 例 3-30 PHP 中的引用

```
//正则表达式只包含 9 个捕获分组，将捕获的文本替换为空字符串
echo preg_replace("/^(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)/", "$10", "0123456789");
9

//正则表达式包含 10 个捕获分组，将捕获的文本替换为 10 号分组匹配的 9
echo preg_replace("/^(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)/", "$10",
"0123456789");
9

Echo preg_replace("/^(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)/", "${1}0",
"0123456789");
00
```

注：正则表达式两端的/是分隔符，PHP 规定正则表达式两端必须使用分隔符。

Python 和 PHP 的规定明确，所以避免了 `\num` 的二义性；其他一些语言却不是如此，根据它们的文档，引用捕获分组只有 `\num`（或者 `$num`）一种记法，这时候 `\10`（其实 `\11`、`\21` 等都是如此）的二义性问题就无可避免了（实际上，本书中介绍的语言，除了 Python 和 PHP 之外都是如此）。

比如 Java 对 `\num` 中的 `num` 是这样规定的：如果是一位数，则引用对应的捕获分组；如果是两位数且存在对应捕获分组时，引用对应的捕获分组，如果不存在对应的捕获分组，则引用一位数编号的捕获分组。

也就是说，如果确实存在编号为 10 的捕获分组，则 `\10` 引用此捕获分组匹配的文本；否则，`\10` 表示“第 1 个捕获分组匹配的文本”和“字符 0”。程序的运行结果见例 3-31。

## 例 3-31 Java 中的引用

```
//存在 10 分组
System.out.println("0123456789".replaceAll("^(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)$", "$10"));
9

//不存在 10 分组
System.out.println("012345678".replaceAll("^(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)(\d)$", "$10"));
00
```

除 Java 之外，Ruby 和 JavaScript 也采用这种规定，它看起来有点古怪，而且有个问题无法解决：如果存在编号为 10 的捕获分组，无法用 `\10` 表示“编号为 1 的捕获分组和字符 0”，因为此时 `\10` 表示的必然是编号为 10 的捕获分组。

在开发中，尤其是进行文本替换时有时确实会遇到这个问题，在现有的规则下是无解的。好

在，一般我们并不会用到太多的捕获分组（包含捕获分组数超过 10 个的表达式很少见，也很难理解和维护）。而且，已经有越来越多的语言提供了命名分组，它可以彻底解决这个问题。

3.3.3 命名分组

捕获分组通常用数字编号来标识，但这样有几个问题：数字编号不够直观，虽然规则是“从左向右按照开括号出现的顺序计数”，但括号多了难免混淆；引用时也不够方便，上面已经讲过 \10 引起混淆的情况。

为解决这类问题，一些语言和工具提供了**命名分组**（named grouping），可以将它看做另一种捕获分组，但是标识是容易记忆和辨别的名字，而不是数字编号。

命名分组的记法也并不复杂。在 Python 中用 (?P<name>...) 来分组的，其中的 name 是赋予这个分组的名字，regex 则是分组内的正则表达式。这样，匹配年月日的正则表达式中，可以给年、月、日的分组分别命名，再用 group(name) 来获得对应分组匹配的文本。图 3-5 说明了命名分组的结构，具体的代码见例 3-32。

字符串:        2010        -        12        -        22  
字符串:        2011        -        01        -        03  
  
表达式: (?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})  
分组名:        year                    month                    dat

图 3-5 命名分组

例 3-32 命名分组捕获

```
namedRegex = r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})"  
result = re.search(namedRegex, "2010-12-22")  
print result.group("year")  
2010  
  
print result.group("month")  
12  
  
print result.group("day")  
22
```

因为数字编号分组的历史更长，为保证向后兼容性，即便使用了命名分组，每个命名分组同时也具有数字编号，其编号规则没有变化。从例 3-33 可以看到，在全部使用命名分组的情况下，仍然可以使用数字编号来引用分组。

例 3-33 命名分组捕获时仍然保留了数字编号

```
namedRegex = r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})"  
result = re.search(namedRegex, "2010-12-22")
```

```
print result.group(1)
2010

print result.group(2)
12

print result.group(3)
22
```

在 Python 中，如果使用了命名分组，在表达式中反向引用时，必须使用(*?P=name*)的记法；而要进行正则表达式替换，则需要写作\g<name>，其中的 *name* 是分组的名字。代码见例 3-34。

例 3-34 命名分组的引用方法

```
re.search(r"^(?P<char>[a-z])(?P=char)$", "aa") != None # => True

re.sub("(?P<digit>\d)", r"\g<digit>0", "123");
102030
```

值得注意的是，命名分组不是目前通行的功能，不同语言的记法也不同，表 3-3 总结了目前常见的用法。

表 3-3 不同语言中命名分组的记法

语言	分组记法	表达式中的引用记法	替换时的引用记法
.NET	(?<name>...)	\k<name>	\${name}
PHP	(?P<name>...)	(?P=name) <sup>1</sup>	不支持，只能使用\\${num}，其中 num 为对应分组的数字编号
Python	(?P<name>...)	(?P=name)	\g<name>
Ruby	(?<name>...)	\k<name>	\k<name>

注 1：Java 5 和 Java 6 都不支持命名分组，根据目前看到的 JRE 的文档，Java 7 开始支持命名分组<sup>2</sup>，其记法与 .NET 相同。

注 2：Ruby 1.9 以上版本才支持使用命名分组。

### 3.4 非捕获分组

目前为止，总共介绍了括号的三种用途：**分组**，将相关的元素归拢到一起，构成单个元素；**多选结构**，规定可能出现的多个子表达式；**引用分组**，将子表达式匹配的文本存储起来，供之后引用。

<sup>1</sup> 在 PHP 5.2.2 以后可以使用\k<name>或者\k'name'，在 PHP 5.2.4 之后可以使用\k{name}和\g{name}。

<sup>2</sup> 可参考 <http://cr.openjdk.java.net/~sherman/6350801/webrev.00/regex/Pattern.html#groupname>。

这三种用途并不是彼此独立的，而是互相重叠的：单纯的分组可以视为“只包含一个多选分支的多选结构”；整个多选结构也会被视为单个元素，可以由单个量词限定。最重要的是，无论是否需要引用分组，只要出现了括号，正则表达式在匹配时就会把括号内的子表达式存储起来，提供引用。如果并不需要引用，保存这些信息无疑会影响正则表达式的性能；如果表达式比较复杂，要处理的文本又很多，更可能严重影响性能。

为解决这种问题，正则表达式提供了**非捕获分组**（non-capturing group），非捕获分组类似普通的捕获分组，只是在开括号后紧跟一个问号和冒号（?:...），这样的括号叫做非捕获型括号，它只能限定量词的作用范围，不捕获任何文本。在引用分组时，分组的编号同样会按开括号出现的顺序从左到右递增，只是必须以捕获分组为准，非捕获分组会略过，如例 3-35 所示。

### 例 3-35 非捕获分组的使用

```
print re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(2)
12

print re.search(r"(?:\d{4})-(\d{2})-(\d{2})", "2010-12-22").group(1)
12
```

非捕获分组不需要保存匹配的文本，整个表达式的效率也因此提高，但是看起来不如捕获分组美观，所以很多人不习惯这种记法。不过，如果只需要使用括号的分组或者多选结构的功能，而没有用到引用分组，则应当尽量使用非捕获型括号。

如果不习惯这种记法，比较好的办法是，在写正则表达式时统一使用捕获分组，确保正确之后，再把不需要引用的分组修改为非捕获分组——当然，引用分组的编号可能也要调整（上例中，只需要取月份信息，把第一个分组改为非捕获分组之后，取月份信息对应分组的编号从 2 变为 1）。

在本书中，为了使代码简洁和易于，除非特殊标注，否则不管匹配完成之后是否会引用文本，都使用捕获分组。

## 3.5 补充

### 3.5.1 转义

之前讲到，如果元字符是单个出现的，直接添加反斜线字符转义即可转义，所以\*、+、?的转义形式分别是\\*、\+、\?。如果元字符是成对出现的，则有可能只对第一个字符转义，比如{6}和[a-z]的转义分别是\{6}和\[a-z]。

括号的转义与它们都不同，与括号有关的所有三个元字符(、)、|都必须转义。因为括号非常重要，所以无论是开括号还是闭括号，只要出现，正则表达式就会尝试寻找整个括号，如果只

转义了开括号而没有转义闭括号，一般会报告“括号不匹配”的错误。另一方面，多选结构中的 `|` 也必须转义（多选结构可以不用括号只出现 `|`），所以，也不要忘记对 `|` 的转义；否则就可能出现例 3-36 的问题。

### 例 3-36 括号的转义

```
re.search(r"^\(a\)$", "(a)") != None      # => True
re.search(r"^\(a\)$", "(a)") != None      # => True
re.search(r"^\(a\)$", "(a)") != None      # => True
Traceback (most recent call last):
error: unbalanced parenthesis

#未转义|
re.search(r"^\(a|b\)$", "(a|b)") != None   # => False
#同时转义了|
re.search(r"^\(a|b\)$", "(a|b)") != None   # => True
```

## 3.5.2 URL Rewrite

提到括号的分组和引用功能，就不能不提到 URL Rewrite。URL Rewrite 是常见 Web 服务器中都具备（也必须）的功能，它用来进行网址的转发，下面是一个转发的例子。

```
外部访问 URL
http://www.example.com/blog/2006/12
内部实现
http://www.example.com/blog/posts.php?year=2006&month=12
```

这样的好处是隔离了外部接口和内部实现，方便修改；也有利于提供更有意义、更直观的 URL。

一般来说，URL Rewrite 都是使用转发规则实现的，每条转发规则对应一类 URL，以正则表达式解析并提取出所需要的信息，重组之后再转发。比如上面的转发，就需要先提取年、月、日的信息进行重组。很自然地，我们会想到使用括号和引用分组的功能来实现。下面就以刚才提到的日期转发为例，看上面的转发规则在当前主流的 Web 服务器中如何配置。

### Microsoft IIS

在 Web.config 配置文件中，找到 `<rewrite>` 节点，在 `<rules>` 下新增下面的代码。

```
<rule name="Rewrite Rule">
<match url="^blog/([0-9]{4})/([0-9]{2})/?$" />
<action type="Rewrite" url="blog/posts.php?year={R:1}&month={R:2}" />
</rule>
```

其中<match>节点中的 url 是外部访问的 URL。对转发的 URL 而言,能接收的都是 path 部分,如果 URL 是 `http://www.example.com/blog/2006/12`,则 path 就是 `blog/2006/12`。正则表达式以 `^blog` 开头,分别用 `[0-9]{4}`、`[0-9]{2}` 匹配其中的年、月信息,因为之后的转发需要用到这些信息,所以必须使用捕获分组以便引用。另外,因为 URL 最后可能出现反斜线/,也可能不出现,意义没有区别,所以使用了量词 `/?`。

Action 节点中的 url 则是转发之后(也就是内部使用)的 URL,转发到 `blog/posts.php`,且将年、月信息作为请求参数,附在后面。在 IIS 中,通过 `{R:num}` 的记法引用分组,其中 `num` 为对应分组的编号;另外,因为这是一个 XML 文件,所有的 `&` 必须转义为 `&amp;`,URL 中的 `&` 也不例外。

关于 IIS 中 URL Rewrite 的具体信息,可以参考下面的详细文档。

<http://learn.iis.net/page.aspx/496/iis-url-rewriting-and-aspnet-routing/>

### Apache

在 `httpd.conf` 配置文件中,找到虚拟主机对应的配置字段,首先确认启用了 URL Rewrite 功能,也就是保证出现了下面这行:

```
RewriteEngine on
```

然后编写规则,上面的转义对应的规则如下:

```
RewriteRule ^blog/([0-9]{4})/([0-9]{2})/?$ blog/posts.php?year=$1&month=$2 [L]
```

以 `RewriteRule` 开头的行指定了转发规则,`RewriteRule` 之后是外部 URL 和转发的 URL,最后是可选出现的标志位(flags,[L]表示“如果 URL 匹配成功,按本条规则转发之后,不再考虑其他转发规则”),这几个字段之间用任意空白字符分隔。在 Apache 中,分组的引用使用 `$num` 的形式,其中 `num` 为分组对应的编号。

关于 Apache 中 URL Rewrite 的具体信息,可以参考下面的详细文档。

Apache 2.x 版: <http://httpd.apache.org/docs/2.0/misc/rewriteguide.html>

Apache 1.3 版: [http://httpd.apache.org/docs/1.3/mod/mod\\_rewrite.html](http://httpd.apache.org/docs/1.3/mod/mod_rewrite.html)

### Nginx

在 `Nginx.conf` 配置文件中找到对应虚拟主机的配置字段,在其中添加下面的规则。

```
rewrite ^blog/([0-9]{4})/([0-9]{2})/?$ blog/posts.php?year=$1&month=$2 last;
```

以 `rewrite` 开头的行指定了转发规则,`rewrite` 之后是外部 URL 和转发的 URL,最后是可选出现的标志位(flags,last 的含义与 Apache 转发规则中的[L]相同),这几个字段之间也是用任意空白字符分隔(要注意,行的末尾必须有分号;)。在 Nginx 中,使用 `$num` 的记法引用分组,

其中 *num* 为分组对应的编号。

相对来说，Nginx 的转发功能最为强大，因为 Apache 和 IIS 的转发一般都只限于单条语句，但是 Nginx 的转发可以使用复杂的判断逻辑，比如下面的转发首先判断浏览器的 user-agent，如果是 IE 则转发，否则不转发。

```
if ($http_user_agent ~ MSIE) {
rewrite ^blog/([0-9]{4})/([0-9]{2})/?$ blog/posts.php?year=$1&month=$2 last;
}
```

关于 Nginx 中 URL Rewrite 的具体信息，可以参考下面的详细文档。

<http://wiki.nginx.org/HttpRewriteModule>

### 3.5.3 一个例子

这部分内容来自一位朋友的问题，这个问题相当有迷惑性和代表性，所以不妨列在这里，希望能解开更多读者的类似疑惑。

问题是这样的：运行 `re.findall('(\w+\.?)+', 'aaa.bbb.ccc')`，期望得到序列 `aaa.`、`bbb.`、`ccc`，实际运行的结果却只有 `ccc`，这是为什么呢？

其实答案很简单——因为表达式 `(\w+\.?)+` 中存在量词 `+`，所以整个正则表达式的匹配过程中，括号内的 `\w+\.?` 会多次匹配：第 1 次匹配 `aaa.`，第 2 次匹配 `bbb.`，第 3 次（也就是最后）匹配 `ccc`，最终这个捕获分组匹配的文本就是 `ccc`。调用 `re.findall()` 时，因为存在括号（也就是捕获分组），默认返回捕获分组匹配的文本，也就是 `ccc`。

解答了这个问题之后，他继续问：如果字符串是 `aaa.bbb`，或者 `aaa.bbb.ccc.ddd`，如何能用一个表达式，逐个拆分出 `aaa.`、`bbb.` 之类的子串呢？（请注意，子串的个数是变化的，并且不能预先知道。）

要解答这个问题，需要记住：捕获分组的个数是不能动态变化的——单个正则表达式里有多少个捕获分组，一次匹配成功之后，结果中就必然存在多少个对应的元素（捕获分组匹配的文本），如果不能预先规定匹配结果中元素的个数，就不能使用捕获分组。如果要匹配数目不定的多段文本，必须通过重复多次匹配完成。具体到这个例子，在 `re.findall('(\w+\.?)+', 'aaa.bbb.ccc')` 中，整个正则表达式会匹配成功 3 次，得到 3 个子串；如果把字符串改为 `aaa.bbb.ccc.ddd`，则整个正则表达式会匹配成功 4 次，得到 4 个子串。

# 第 4 章 断言

正则表达式中的大多数结构匹配的文本会出现在最终的匹配结果中(一般用 `group(0)`可以得到),但是也有些结构并不真正匹配文本,而只负责判断在某个位置左/右侧的文本是否符合要求,这种结构被称为**断言**(`assertion`)。常见的断言有三类:单词边界、行起始/结束位置、环视。

## 4.1 单词边界

在文本处理中经常可能进行单词替换,比如把一段文本中的 `row` 都替换成 `line`。一般想到的是调用字符串的替换方法,直接替换 `row`。在不同语言中这些方法各不相同,但差别不大。

```
替换前: The row we are looking for is row 10.  
替换后: The line we are looking for is line 10.
```

不过,这样替换也可能会造成意想不到的后果。

```
替换前: ...tomorrow I will wear in brown standing in row 10 next to the rowdy guy...  
替换后: ...tomorline I will wear in blinen standing in line 10 next to the linedy guy...
```

不仅所有单词 `row` 都被替换成了 `line`,其他单词内部的 `row` 也被替换成了 `line`,这显然不是我们想要的结果。

要解决这个问题,必须有办法确定**单词 row**,而不是**字符串 row**。为解决这类问题,正则表达式提供了专用的**单词边界**(`word boundary`),记为 `\b`。它匹配的是“单词边界”位置,而不是字符。也就是说, `\b` 能够匹配这样的位置:一边是单词字符,另一边不是单词字符,如图 4-1 所示。

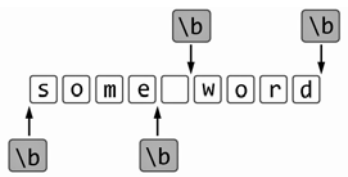


图 4-1 单词边界的匹配

表 4-1 更详细地说明了 `row` 配合单词边界 `\b` 之后的匹配情况。



表 4-1 单词边界的匹配

字符串	<code>\brow\b</code>	<code>\brow</code>	<code>row\b</code>
tomorrow			✓
brown			
row	✓	✓	✓
rowdy		✓	
表达式说明	只能是单词 row	<code>\b</code> 的右侧是单词字符，所以左侧不能是单词字符	<code>\b</code> 的左侧是单词字符，所以右侧不能是单词字符

观察表格，可以发现两点：第一，单词边界并不区分左右，在“单词边界”上，可能只有左侧是单词字符，也可能只有右侧是单词字符，总的来说，单词字符只能出现在一侧；第二，单词字符要求“另一边不是单词字符”，而不是“另一边的字符不是单词字符”，也就是说，一边必须出现单词字符，另一边可以出现非单词字符，也可能没有任何字符。所以，如果字符串只包含单词 `word`，用 `\bword\b` 应该是可以匹配的，虽然 `w` 之前和 `d` 之后都没有任何字符。

单词边界要求一侧必须出现单词字符，到底什么是单词字符呢？

一般情况下，“单词字符”的解释是 `\w` 能匹配的字符。在 JavaScript、PHP、Python 2、Ruby 中，`\w` 只能匹配 `[0-9a-zA-Z_]`。所以在这些语言中，`\b\w+\b` 能准确匹配英文单词了。给定一段文本，就可以用 `\b\w+\b` 将所有的单词提取出来，例如 4-1。

例 4-1 单词边界的匹配

```
print re.findall(r"\b\w+\b", "a sentence\tcontains\na lot of words")
['a', 'sentence', 'contains', 'a', 'lot', 'of', 'words']
```

在 Web 开发中，经常需要对某些单词标记高亮，一般的做法是在单词的前后加上 tag，比如 `<span class="hl">` 和 `</span>`，这个功能也可以由单词边界配合正则表达式替换完成，代码见例 4-2。

例 4-2 使用单词边界，准确给单词标记高亮

```
string = "tomorrow I will wear in brown standing in row 10 next to the rowdy guy"
rowWordRegex = r"\brow\b"
print re.sub(rowWordRegex, r"<span class='hl'>\g<0></span>", string)
tomorrow I will wear a brown sweater standing in row <span class="hl">10</span>
next to the rowdy guy
```

但是也有些单词，`\b\w+\b` 是无能为力的，比如 `e-mail` 和 `M.I.T.`。因为连字符和点号都不能由 `\w` 匹配，所以 `\b\w+\b` 无法匹配 `e-mail`，也无法匹配 `M.I.T.`。如果确实希望处理 `e-mail`

之类的“单词”，也可以把表达式改为**`\b[-\w]+\b`**。

如果使用的是.NET 或者 Python 3，在默认情况下，**`\w`** 不仅仅等价于 **`[0-9a-zA-Z_]`**，还能匹配各种语言中的“单词字符”，包括中文字符，这时**`\b`**的使用就有很大不同，具体情况（[☞59](#)）现在不展开。

与单词边界**`\b`**对应的还有非单词边界**`\B`**，两者的关系类似**`\s`**和**`\S`**、**`\w`**和**`\W`**、**`\d`**和**`\D`**：在同一种语言中<sup>1</sup>，不管**`\b`**是如何规定的，**`\b`**能匹配的位置，**`\B`**就不能匹配；**`\B`**能匹配的位置，**`\b`**就不能匹配。但是在实际使用中，**`\B`**使用频率远远少于**`\b`**，所以暂不做过多介绍。

## 4.2 行起始/结束位置

单词边界匹配的是某个位置而不是文本，在正则表达式中，这类匹配位置的元素叫做**锚点**（anchor），它用来“定位”到某个位置。常用的锚点还有**`^`**和**`$`**，它们分别匹配字符串的开始位置和结束位置，所以可以用来判断“整个字符串能否由表达式匹配”。这两个锚点在第 1 章粗略介绍过，现在仔细讲解。

一般情况<sup>2</sup>下，**`^`**匹配**整个字符串的起始位置**。

字符串	Some sample text
<b><code>^</code></b> 能匹配的位置	↑

依靠**`^`**，就可以用正则表达式**`^Some`**准确验证字符串“是否以 **Some** 开头”，因为**`^`**会把整个表达式的匹配“定位”在字符串的开始位置。这样，即便表达式的其他部分可以在字符串中其他位置找到匹配，整个表达式也无法匹配成功。

在某些情况下，**`^`**也可以匹配字符串内部的“行起始位置”。在讲解这种情况之前，我们先来看看怎么划分行。

在编辑文本时，敲回车键就输入**行终止符**（Line terminal），结束当前行，新起一行。看起来，这很好理解，然而不同平台上的行终止符其实各不相同，表 4-2 列出了常见平台下的“行终止符”（这里不说“回车”而说“行终止符”是为了避免混淆，因为**`\r`** 字符就叫“回车（符）”，而**`\n`** 字符叫做“换行符”，“这里的回车是由回车符和换行符构成的”之类说法容易引起误解）。

<sup>1</sup> 更准确的说法是，“在同一种语言的同一种匹配模式下”，因为在.NET 或 Python 中，可以设定匹配模式，指定**`\w`** 采用不同的匹配规则，继而影响**`\b`** 的匹配；另外 locale（语言环境）的设定也有可能影响到**`\w`** 和**`\b`** 的匹配，只是这种情况比较少见，所以此处不介绍。

<sup>2</sup> 这里说的“一般情况”是指在本书提到的各编程语言中默认情况，只有 Ruby 是例外，下文会详细介绍。

表 4-2 不同平台下的行终止符

平台	行终止符
UNIX/Linux	\n
Windows	\r\n
Mac OS <sup>1</sup>	\n

也就是说，每一行的“起始位置”，就是“行终止符”之后的那个位置，如果没有专门的符号，就要考虑各种“行终止符”。下面的例子看得更清楚，为了让换行符“可见”，我们用 `NL` 表示。

```
first line
second line
last line
```

它其实是下面这样，其中的 `NL` 可能是 `\n`，也可能是 `\r\n`。

```
first lineNLmiddle lineNLlast lineNL
```

如果把匹配模式设定为**多行模式**（Multiline Mode，这是一种影响元字符匹配的设定，下一章详细讲解）下，`^` 就既可以匹配整个字符串的起始位置，也可以匹配换行符之后的位置（设定多行模式最简单的办法是在正则表达式之前加上 `(?m)`，这里虽然出现了括号，但因为是专用于指定匹配模式，所以不会作为捕获分组）。

表达式	(?m)^
字符串	first lineNLmiddle lineNLlast lineNL
^能匹配的位置	↑                    ↑                    ↑                    ↑

注意，如果字符串的末尾出现了行终止符，`^` 也会匹配这个行终止符之后的位置。这样做是有意义的，比如要在每行开头加上特殊标注（最常见的处理是把纯文本格式转换为 HTML 格式），末尾的空行自然不应该漏掉。不过一般来说，`^` 的主要用途是与其他子表达式配合，比如例 4-3 那样，提取每行的第一个单词。

例 4-3 提取每行的第一个单词

```
string = "first line\nsecond line\r\n\r\nlast line"
lineBeginWordRegex = r"(?m)^\w+"
print re.findall(lineBeginWordRegex, string)
['first', 'second', 'last']
```

如果不想定位到字符串内部的行起始位置，只关心整个字符串的起始位置，则可以使用 `\A`，

<sup>1</sup> 有些资料上说 Mac OS 中的换行符是 `\r`，这并不算错，但只有 OS 9 之前的版本中，换行符才是 `\r`。从 OS X 开始，换行符与 UNIX/Linux 一样都是 `\n`，考虑到 OS 9 距今已经很久了，本书不考虑 `\r` 的情况。参考 <http://en.wikipedia.org/wiki/Newline>。

绝大多数工具中的正则表达式都支持这个锚点，它在任何情况下（包括多行模式下）都只匹配整个字符串的起始位置，如例 4-4 所示。

例 4-4 匹配整段文本的第一个单词

```
string = "first line\nsecond line\r\nast line"
stringBeginWordRegex = r"(?m)\A\w+"
print re.findall(stringBeginWordRegex, string)
['first']
```

“行结束位置”的情况更复杂。除去“行终止符”可能由各种字符表示的情况之外，“行结束位置”可能没有任何字符，还是上面的字符串，你猜猜它有几个行终止符？

字符串: Some sample text

可能一: Some sample text

可能二: Some sample textNL

而且其中的 NL 可能是\n，也可能是\r\n。

如果要匹配字符串最后一个单词，不但必须考虑 NL 所对应字符的多种可能，而且要兼顾 NL 是否出现情况更加复杂。

针对这种问题，正则表达式提供了“通吃”行结束符的锚点\$，它匹配的同样是位置。通常它匹配的是整个字符串的结尾位置——如果最后是行终止符，则匹配行终止符之前的位置；否则，匹配最后一个字符之后的位置。

字符串	Some sample text
\$能匹配的位置	↑
字符串	Some sample textNL
\$能匹配的位置	↑

这时候，无论 NL 是什么，是否存在，都可以用表达式 \w+\$ 匹配最后一个单词，代码见例 4-5。

例 4-5 匹配整段文本的最后一个单词

```
stringEndWordRegex=r"\w+$"
string = "Some sample text"
print re.findall(stringEndWordRegex, string)
['text']

string="Some sample text\n"
print re.findall(stringEndWordRegex, string)
['text']
```

如果指定了多行模式，`$`会匹配每个行终止符之前的位置。如果最后一行没有行终止符，则匹配字符串的结尾位置；否则，匹配行终止符之前的位置。

表达式	<code>(?m)\$</code>
字符串	first lineNLmiddle lineNLlast line
<code>\$</code> 能匹配的位置	↑ ↑ ↑
字符串	first lineNLmiddle lineNLlast lineNL
<code>\$</code> 能匹配的位置	↑ ↑ ↑

如果指定了多行模式，就可以用`\w+$`匹配每一行的最后一个单词了，代码见例 4-6。

#### 例 4-6 匹配每行的最后一个单词

```
string = "first line\nsecond line\r\nlast line"
lineBeginWordRegex = r"(?m)\w+$"
print re.findall(lineBeginWordRegex, string)
['line', 'line', 'line']
```

与`$`类似的还有两个特殊标记`\Z`和`\z`，它们不受多行模式的影响，在任何情况下都匹配整个字符串的结束位置。`\Z`和`\z`的主要差别在于：`\Z`等价于默认模式（非多行模式）下的`$`，如果字符串的末尾有行终止符，则它匹配换行符之前的位置；`\z`则不管行终止符，只匹配“整个字符串的结束位置”。

字符串	first lineNLmiddle lineNLlast line
<code>\Z</code> 能匹配的位置	↑
<code>\z</code> 能匹配的位置	↑
字符串	first lineNLmiddle lineNLlast lineNL
<code>\Z</code> 能匹配的位置	↑
<code>\z</code> 能匹配的位置	↑

回过头来说`^`和`$`，这两个锚点第一章介绍过，从数据校验的例子可以看到，`re.search(pattern, string)`只表示 `pattern` 能否在 `string` 中找到匹配，但是如果 `pattern` 只匹配 `string` 中的一部分，也不会返回 `None`；为了验证整个 `string` 能否由 `pattern` 匹配，通常的做法是在 `pattern` 两端加上`^`和`$`，就像例 4-7 那样。

#### 例 4-7 借助`^`和`$`完成数据验证

```
#没有使用^和$, 只匹配子串 123456
re.search(r"\d{6}", "a123456b") != None      # => True
#使用^和$, 完整验证
re.search(r"^\d{6}$", "a123456b") != None   # => False
re.search(r"^\d{6}$", "123456") != None     # => True
```

常用到数据验证的场合就是对用户提交的数据进行验证，比如在网页上，要求用户在输入框（比如密码、邮箱）中填入某些信息，加以验证。一般来说，输入框中都不能输入换行符，但如果用户使用程序来提交，接收到的值就可能包含换行符。在这种情况下，在正则表达式两端添加`^`和`$`是无法准确验证的，因为`$`可以匹配“结尾行终止符之前的位置”，验证时就忽略了末尾的行终止符。使用`\z`替换`$`可以堵住这个漏洞（使用`\z`时，最好把`^`也替换成`\A`，这样更符合习惯，因为一般`^`是和`$`联系的）。下面用 Python 为例说明这一点（因为 Python 不支持`\z`，但是 Python 中的`\Z`等价于其他语言中的`\z`，所以例 4-8 使用`\Z`）。

#### 例 4-8 借助`\A`和`\Z`完成更准确的数据验证

```
#只使用^和$验证 6 位数字字符串，可能有漏洞
re.search(r"^\d{6}$", "123456") != None      # => True
re.search(r"^\d{6}$", "123456\n") != None    # => True
#Python 不支持\z，它的\Z 等价于其他语言中的\z
re.search(r"\A\d{6}\Z", "123456") != None     # => True
re.search(r"\A\d{6}\Z", "123456\n") != None   # => False
```

类似 Python 中的`\Z`，JavaScript 中的`$`的匹配也比较特殊，JavaScript 没有提供`\A`、`\z`、`\Z`，只有`^`和`$`，但是 JavaScript 中的`$`，只能匹配字符串/行的结束位置，即便字符串末尾有换行符，也是如此，所以在验证时，可以放心使用`^`和`$`。JavaScript 代码见例 4-9。

#### 例 4-9 JavaScript 中的验证

```
/^a$/.test("a\n");      // => False
/^a$/.test("a");         // => True
```

`^`和`$`的另一个特点是，进行正则表达式替换时并不会被替换。也就是说，在起始/结束位置进行替换，只会在起始/结束位置添加一些字符，位置本身仍然存在。使用这个特性，我们可以很方便地转换文本的格式。常见的应用是将纯文本转换为 HTML，比如将纯文本的电子文档转换成 epub 格式，就需要如此处理。这个问题最简单的思路是，使用多行模式将`^`替换为`<p>`，将`$`替换为`</p>`，如例 4-10 所示，注意其中使用了多行模式，这样可以找到字符串内部文本行的开始和结束位置。

#### 例 4-10 `^`和`$`的替换

```
plainText = "line1\nline2\nline3"
print plainText
line1
line2
line3

print re.sub(r"(?m)$", "</p>", re.sub(r"(?m)^", "<p>", plainText))
<p>line1</p>
<p>line2</p>
<p>line3</p>
```

`^`和`$`的另一个常用功能是删去多余的空白，包括行首尾的空白和空行。因为种种原因，要处理的文本可能经常包含许多空白字符，有些出现在行首，有些出现在行尾，还可能有不少空行。比如下面这段文本（为方便识别，在行的首尾分别用「和」标识）。

```
「  begin」
「  between 」
「」
「end   」
```

如果要整理格式，需要删掉不必要的空白字符，但又不能把所有空白字符都删掉（单词与单词之间的空白字符应当保留）。所以要做的其实是删除行首和行尾的空白字符，我们先删除行首的空白字符，使用的正则表达式是`(?m)^\s+`。

这里必须使用多行模式，否则就只能删除整个字符串首尾的空白字符。另一方面，此处使用了量词`+`而不是`*`，因为`^\s*`可以不匹配任何字符，这样的“删除”没有意义。将`(?m)^\s+`匹配的文本替换为空字符串，就执行了删除操作（一般正则表达式应用中没有单独的“删除”操作，删除操作都是通过将文本替换为空字符串实现的）<sup>1</sup>。例 4-11 所示展示了去除字符串行首空白字符的代码。

#### 例 4-11 去除行首的空白字符

```
withSpaces = "  begin\n between\t\n\nend  "
beginSpaceRegex = r"(?m)^\s+"
trimmedLeadingSpace = re.sub(beginSpaceRegex, "", withSpaces)
print trimmedLeadingSpace
```

为方便观察，仍然用上面的方式显示字符串。

```
「begin」
「between 」
「end   」
```

因为`\s`匹配的空白字符中包含换行符`\n`，所以完全是空白字符的第 3 行（连同换行符）也被删掉了。现在来删行尾的空格，使用表达式`\s+$`，同样要记得使用多行模式，如例 4-12 所示。

#### 例 4-12 去除行尾的空白字符

```
#代码接上面
endSpaceRegex = r"(?m)\s+$"
```

---

<sup>1</sup> 在实际开发中，经常可能要进行这种处理，比如 Web 页编辑器为了显示缩进，会在行的开头加入很多空白字符，这些空白字符都会消耗流量，所以正式发布时，可以用`(?m)^\s+`全部删除。

```
trimmedEndingSpace = re.sub(endSpaceRegex, "", trimmedLeadingSpace)
print trimmedEndingSpace
```

仍然用上面的方式显示字符串。

```
「begin」
「between」
「end」
```

能不能用多选结构 `(^\s+|\s+$)` 并列两个表达式，一步完成呢？

```
withSpaces = "    begin\n between\t\n\nend  "
spaceRegex = r"^(?m)(^\s+|\s+$)"
spaceTrimmed = re.sub(spaceRegex, "", withSpaces)
print spaceTrimmed

「begin」
「betweenend」
```

不但第三行被删除，第二行和第四行也合并成一行，中间的 `\t\n\n` 全部删除了，第二行末尾没有了换行符；而真正的目的其实只是将 `\t\n\n` 替换为 `\n`。仔细看看正则表达式 `(^\s+|\s+$)` 可以知道，在 `\s+$` 中，`\s` 可以匹配 `\t` 和 `\n`，所以 `\s+$` 可以匹配开始的 `\t\n`，同样 `^\s+` 可以匹配结尾的 `\n`，所以 `\t\n\n` 经过两步被彻底删除了。

这个例子所用的表达式很有意思<sup>1</sup>，它提醒我们用多选结构合并多个表达式时，一定要小心未曾预期的后果；有时候，分几步进行反而能省去许多麻烦，这类例子在第 9 章还要讲到。

下面的表格总结了各语言中 `^`、`$`、`\z` 和 `\Z` 的匹配情况。其中 Ruby 是例外，Ruby 提供了多行模式，但这个“多行模式”其实等价于常说的“单行模式”，它的作用只是让点号 `.` 能够匹配合换行符（下一章会详细讲解各种模式），而不影响 `^` 和 `$` 的匹配；另一方面，Ruby 默认模式已经是“多行模式”，所以 `$` 可以匹配字符串内部的行结束位置。

表 4-3 `^`和`$`的总结

模式	行为	.NET	Java	JavaScript	PHP	Python	Ruby
默认模式	<code>^</code> 匹配字符串起始位置	✓	✓	✓	✓	✓	✓
	<code>^</code> 匹配字符串内部行起始位置						✓
	<code>\$</code> 匹配字符串结束位置	✓	✓	✓	✓	✓	✓
	<code>\$</code> 匹配字符串末尾行终止符之前	✓	✓		✓	✓	✓
	<code>\$</code> 匹配字符串内部行结束位置						✓

<sup>1</sup> 实际上，《精通正则表达式》（第 3 版）也详细讲到了这个例子，有兴趣的读者可以参考该书第 199 页。



( 续表 )

模式	行为	.NET	Java	JavaScript	PHP	Python	Ruby
支持多行模式		√	√	√	√	√	
多 行 模 式	^匹配字符串起始位置	√	√	√	√	√	无此模式
	^匹配字符串内部行起始位置	√	√	√	√	√	无此模式
	\$匹配字符串结束位置	√	√	√	√	√	无此模式
	\$匹配字符串内部行结束位置	√	√	√	√	√	无此模式
	\A 等于默认模式的^	√	√		√	只能匹配字符串的起始位置	×
	\Z 等于默认模式的\$	√	√		√	只能匹配字符串的结束位置	×
	\z 匹配字符串的结束位置	√	√		√	无	×

注 1：JavaScript 只支持^和\$，而且\$不会匹配末尾换行符之前的位置。

注 2：Ruby 默认采用了多行模式，所以其中\A、\Z、\z 能匹配的位置与其他语言中的相同，但 Ruby 中默认模式下的^和\$都可以匹配文本内部行的起始/结束位置，所以\A 不等于默认模式的^，\Z 也不等于默认模式的\$。

4.3 环视

前面介绍过单词边界匹配的是这样的位置：一边是单词字符，另一边不是单词字符。从另一个角度来看，它能进行这样的判断：在某个位置向左/向右看，必须出现或不能出现某类字符。有时候，这种功能非常有用。

在第 2 章，用表达式<[ ^/> ] [ ^> ] \*>匹配 open tag，它保证了<之后不会出现/，这样就排除了</img>之类的 close tag，但它也可以匹配 self-closing tag，比如<br/>。如果将表达式改为<[ ^/ ] [ ^> ] \* [ ^/ ] >，又会有一个问题，在<和>中的[ ^/ ] [ ^> ] \* [ ^/ ]，能匹配的文本至少包含 2 个字符，这样就无法匹配<u>。

上一章用正则表达式<[ ^/ ] ( [ ^> ] \* [ ^/ ] ) ?>解决了这个问题。但是仔细想想，也可以从另一个角度描述：在开始位置匹配<，同时要求这个<之后不能是/；然后匹配中间的文本，除非在属性（也就是引号字符串）中，否则不能出现>，且长度必须大于 1（<>不是合法的 tag）；最后匹配>，同时要求这个>之前不能是/。

这样描述更加准确，逻辑也更清晰，可以用三个子表达式分别匹配这三个部分。单独来看，

开头的<、中间的内容、结尾的>都不难匹配，但必须解决一个问题：匹配开头的<时，除去找到<字符，还必须向后（向右）看看，确认字符不能是/，同时又不能真正匹配这个字符，因为<和>中间的文本是有单独的子表达式匹配的；同样，结尾>的匹配也是如此。

针对这种要求正则表达式专门提供了**环视**（look-around）用来“停在原地，四处张望”。环视类似单词边界，在它旁边的文本需要满足某种条件，而且本身不匹配任何字符。

比如正则表达式<(?!/)，其中的(?!/)是一个环视结构，(?!...)是这个结构的标识，/才是真正的表达式，整个结构的意思是“当前位置之后（右侧），不容许出现/能匹配的文本”。看起来它和<[^/]类似，其实大不相同：如果<(?!/)匹配成功，正则表达式真正匹配完成的只有<，而不包括<之后的那个字符，这样，就能准确表示“匹配<，同时这个<之后不能是/”。

再来看表达式(?<!/)>，其中的(?<!/)>也是一个环视结构，(?<!...)是这个结构的标识，/才是真正的表达式，整个结构的意思是“在当前位置之前（左侧），不容许出现/能匹配的文本”（它与上面的(?!/)类似，只是多了一个<，更加形象地指向左侧）。这样，就能准确地表示“匹配>，同时>之前不能是/”。

至于<和>之间的文本，在上一章曾讲解过，可以用(['^']\*|"[^"]\*"|'['>'])\*+准确匹配。

最后，把这三个部分结合起来，得到正则表达式<(?!/)(['^']\*|"[^"]\*"|'['>'])\*+(?<!/)>，它可以准确匹配 open tag，而且不会错误匹配 self-closing tag，匹配情况如图 4-2 所示，代码见例 4-13。

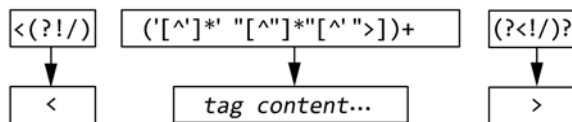


图 4-2 匹配 open tag 的表达式（注意环视结构并不会真正匹配文本）

#### 例 4-13 使用环视结构，准确匹配 open tag

#注意正则表达式内部的双引号字符要转义

```
openTagRegex = r"<(?!/)(['^']*|\"[^\"]*\"|'['>'])*+(?<!/)>Z"
re.search(openTagRegex, "<input name=txt value=\">\">") != None      # => True
re.search(openTagRegex, "<input name=txt value='>'>") != None      # => True
re.search(openTagRegex, "<u>") != None                             # => True
re.search(openTagRegex, "<br/>") != None                             # => False
re.search(openTagRegex, "<img src=\"ur1\"/>") != None               # => False
```

在这个表达式中出现了两种环视：(?!...)和(?<!...)，它们的名字分别是“否定顺序环视”和“否定逆序环视”。“否定”的意思是“如果正则表达式匹配成功，则在当前位置匹配失败”，而“顺序”和“逆序”则表示正则表达式需要匹配的文本所在的位置。所以总的来说，环视一共

分为 4 种：肯定顺序环视（positive-lookahead）、否定顺序环视（negative-lookahead）、肯定逆序环视（positive-lookbehind）、否定逆序环视（negative-lookbehind）。环视的分类见表 4-4。

表 4-4 环视的分类

名字	记法	判断方向	结构内表达式匹配成功的返回值
肯定顺序环视	(?=...)	向右	True
否定顺序环视	(?!...)	向右	False
肯定逆序环视	(?<=...)	向左	True
否定逆序环视	(?<!...)	向左	False

这 4 个名字容易混淆，不妨这样记忆：在当前位置，如果是朝右判断，则是顺序环视（lookahead），如果是朝左判断，则是逆序环视（lookbehind）；如果要求子表达式能匹配的字符串必须出现，则为肯定环视（positive），如果要求子表达式能匹配的字符串不能出现，则为否定环视（negative）。图 4-13 说明，对于字符串 12345，以 \d{3} 为表达式的四种环视能匹配的位置分别是：右侧必须出现三个数字字符，右侧不能出现三个数字字符，左侧必须出现三个数字字符，左侧不能出现三个数字字符。

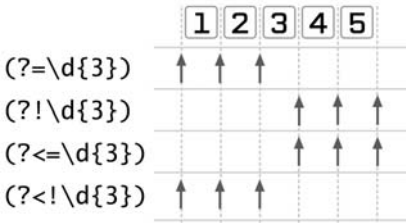


图 4-3 对字符串 12345，4 种环视能匹配的位置

环视的最大特点是“在原地执行”，之前已经看到 <(?!/) 匹配的其实只有一个 < 字符，(?<!/)> 匹配的也只有一个 > 字符，虽然它们都需要测试 / 的匹配。有时候确实需要用到“原地”的判断，因为要寻找的确实只是位置，而不需要真正匹配任何字符，比如整理数字字符串的格式，就是如此。

英文中的数字更习惯用逗号分隔以方便阅读，比如 12345 应该写作 12,345。如果用正则表达式来完成任务，就是“把逗号添加到这样的位置：右侧的数字字符串的长度是 3 的倍数”，看起来只需要使用肯定顺序环视就足够了，用正则表达式找到这样的位置 (?(=(\d{3}))+)，将它“替换”为逗号（其实就是在这里塞进一个逗号），代码见例 4-14。

例 4-14 整理数字字符串，第一次尝试

```
print re.sub(r"(?=(\d{3}))+", ",", "12345")
,1,2,345
```

结果却不是想象的那样。因为“右侧数字字符串”严格说应该是“当前位置右侧，所有数字字符构成的字符串”；但是 `(?=(\d{3})+)` 并不能表达这个意思，比如第一个字符 **1** 之前的位置，右侧数字字符串长度为 5，但其中存在长度为 3 的子串，所以这个位置也可以匹配（如果这样说不够清楚的话，请仔细观察上面的图）。同样，**2**、**3** 之前的位置都是如此。

解决这个问题必须配合否定顺序环视，让 `(\d{3})+` 能匹配右侧的整个数字字符串，而不能只匹配其中的一个子串。也就是说，要一直匹配到“右侧不再有数字字符的位置”为止。所以，必须将表达式改写为 `(?=(\d{3})+(?!\d))`，结果如例 4-15 所示。

#### 例 4-15 整理数字字符串，第二次尝试

```
print re.sub(r"(?=(\d{3})+(?!\d))", r",", "12345")
12,345
```

似乎没问题了，但是从例 4-16 看，如果字符串的长度正好是 3 的倍数，还是有问题。

#### 例 4-16 整理数字字符串，意外的情况

```
print re.sub(r"(?=(\d{3})+(?!\d))", r",", "123456")
,123,456
```

字符串的开头多出了一个逗号，因为这个位置右侧的数字字符串长度为 6。更严格地说，要加入逗号的位置其实是这样的：右侧的数字字符串长度为 3 的倍数，且左侧也是数字字符。所以还需要加上肯定逆序环视，将正则表达式修改为 `(?<=\d)(?=(\d{3})+(?!\d))`，如例 4-17 所示。

#### 例 4-17 整理数字字符串，最后的尝试

```
re.sub(r"(?<=\d)(?=(\d{3})+(?!\d))", r",", "123456")
123,456
```

仔细观察这个例子可以发现，表达式中其实出现了三个环视结构，其中 `(?! \)` 一种组合，除去这个例子中出现的嵌套和并列两种组合，环视结构还可以通过其他方式组合。关于这几类组合，本章末尾有一节专门讲解。

环视是非常有用的功能，日常要执行的许多操作都可能要用到环视，下面再举一个例子。

我们经常会遇到中英文混排的文本，英文文本需要用空白字符来区分单词，中文文本中则很少出现空白字符。但是在转贴或格式转换时，经常会产生一些多余的空白字符：

中 英文混排，some English word，有多余的空 白字符

为了整理格式，需要删除这些空白字符。正则表达式匹配空白字符很容易，直接用 `\s+` 即可。但如果直接删除 `\s+` 能匹配的所有文本，就成了下面这样：

中英文混排, `someEnglishword`, 有多余的空白字符

所以真正要找的, 其实是这样的 `\s+`: 从它向左看, 不能出现英文字母; 从它向右看, 也不能出现英文字母。所以需要在 `\s+` 的两端分别添加否定逆序环视和否定顺序环视, 得到 `(?<![a-zA-Z])\s+(?![a-zA-Z])`, 结果如例 4-18 所示。

#### 例 4-18 去掉中英文混排文本中不必要的空白字符

```
mixedString = "中 英文混排, some English word, 有多余的空 白字符"
spaceBetweenChineseRegex = r"(?<![a-zA-Z])\s+(?![a-zA-Z])"
print re.sub(spaceBetweenChineseRegex, "", mixedString)
中英文混排的文本, some English word, 有多余的空白字符
```

你或许会想, 这个表达式能不能改一改, 比如左侧的否定逆序环视 `(?<![a-zA-Z])`, 能不能改为肯定环视, 指定出现一个非英文字符 `(?<=[^a-zA-Z])`, 右侧的否定顺序环视也改为肯定顺序环视 `(?=[^a-zA-Z])`?

初看起来这并没有问题, 但这个问题其实涉及肯定环视和否定环视的一大根本不同: 肯定环视要判断成功, 字符串中必须有字符由环视结构中的表达式匹配; 而否定环视要判断成功, 却有两种情况: 字符串中出现了字符, 但这些字符不能由环视结构中的表达式匹配; 或者字符串中不再有任何字符, 也就是说, 这个位置是字符串的起始位置或者结束位置。这两种环视的区别, 可以通过例 4-19 展现。

#### 例 4-19 使用不同的环视去掉空白字符

```
#注意字符串首尾的空白字符
mixedString = " 中 英文混排, some English word, 有多余的空 白字符  "
negativeSpaceTrimRegex = r"(?<![a-zA-Z])\s+(?![a-zA-Z])"
#两端加上分隔符便于观察
print "[" + re.sub(spaceBetweenChineseRegex, "", mixedString) + "]"
[中英文混排, some English word, 有多余的空白字符]

positiveSpaceTrimRegex = r"(?<=[^a-zA-Z])\s+(?=[^a-zA-Z])"
print "[" + re.sub(positiveSpaceTrimRegex, "", mixedString) + "]"
[中英文混排的文本, some English word, 有多余的空白字符 ]
```

如果使用肯定环视, 则无法去掉字符串首尾的空白。因为在字符串的开头, `\s+` 虽然能匹配空白字符, 但其左侧并没有任何字符, 所以 `(?<=[^a-zA-Z])` 无法匹配成功; 字符串末尾的 `(?=[^a-zA-Z])` 也是如此。

到现在为止, 如果你觉得自己已经掌握了环视功能, 来看个更复杂的例子: 电子邮件地址中, 更准确的主机名验证。

在上一章，我们给出了一个匹配 E-mail 地址的表达式，但它还不够完整，尤其是主机名部分（hostname）的匹配。根据规范，主机名以点号分隔为多个域名字段（label），每个域名字段可以包含大小写字母、数字字母、横线，但是横线不能出现在最开头位置。关于长度，每个域名字段的长度最多为 63 个字符，整个主机名的长度最多为 255 个字符。通常用的表达式是 `[[-a-zA-Z0-9]{1,63}\.)*[-a-zA-Z0-9]{1,63}`，这个表达式有两个问题：第一，它容许域名字段的第一个字符是横线；第二，它没有限定整个主机名的长度最长为 255 个字符。为准确匹配主机名，就必须解决这两个问题。

为保证域名字段的第一个字符不能是横线，可行的办法之一是单独匹配第一个字符，将表达式 `[[-a-zA-Z0-9]{1,63}\.]` 改写为 `[a-zA-Z0-9][-a-zA-Z0-9]{0,62}\.`；不过，在表达式最开始加上否定顺序环视更加直接，也更加自然，也就是 `(?!-)[-a-zA-Z0-9]{1,63}\.`。

为保证整个主机名字符串长度小于 255 个字符，主机名中全部可能出现的字符都用 `[-a-zA-Z0-9.]` 表示，所以对应的肯定顺序环视就是 `(?=[-a-zA-Z0-9.]{0,255})`，但是并不能直接把它添加到匹配主机名的整个表达式的开头，因为这个表达式只要求匹配一个长度在 255 个字符以内的字符串，并不能保证“之后的整个字符串长度在 255 个字符以内”。如果是单独给出一个字符串，验证它是否是合法的主机名，那么可以在这个环视中的表达式末尾添加 `$`；如果是要从一段文本中提取出某个主机名，那么主机名之后还有其他字符，只是这些字符不能是 `[-a-zA-Z0-9.]`（可能是空白字符，也可能在字符串的末尾，之后没有任何字符），使用否定顺序环视 `(?![-a-zA-Z0-9.])` 就可以兼顾这两种情况。

最终，我们得到了更准确匹配主机名的表达式。

```
(?=[-a-zA-Z0-9.]{0,255})(?![-a-zA-Z0-9.])((?!-)[-a-zA-Z0-9]{1,63}\.)*((?!-)[-a-zA-Z0-9]{1,63})$
```

从例 4-20 可见，这个表达式确实可以更准确地验证主机名。

#### 例 4-20 准确匹配主机名的正则表达式

```
hostnameRegex = r"^(?=[-a-zA-Z0-9.]{0,255})(?![-a-zA-Z0-9.])((?!-)[-a-zA-Z0-9]{1,63}\.)*((?!-)[-a-zA-Z0-9]{1,63})$"
#应该匹配的
re.search(hostnameRegex, "localhost") != None      # => True
re.search(hostnameRegex, "example.com") != None     # => True
#不应该匹配的
#"e" * 64 表示将 64 个"e"连接在一起组成的字符串
re.search(hostnameRegex, "-example.com") != None   # => False
re.search(hostnameRegex, ("e" * 64) + ".com") != None # => False
re.search(hostnameRegex, "e" * 256) != None        # => False
```

## 4.4 补充

### 4.4.1 环视的价值

环视有一个很重要的用途，就是避免编写正则表达式时“牵一发动全身”的尴尬——既可以集中关注某个部分，添加全局性的限制，又不会干扰其他部分的匹配。有些时候，为进行全局性限制而真正匹配文本，会让情况变复杂。

回想匹配 open tag 的例子：open tag 要求<之后不能出现/（否则就是 close tag，比如</a>），而>之前不能出现/（否则就是 self-closing tag，比如<img .../>）；而<和>之间的内容一般使用[<^>]+匹配。如果使用<[<^>][<^>]+[<^>]>，则<和>之间必须出现至少三个字符，即便是<[<^>][<^>]\*[<^>]>，<和>之间也必须至少出现两个字符，都无法匹配<u>，看来很麻烦。

如果使用环视，则非常容易解决：在<之后加上环视(?:!/)，就只匹配<，同时保证匹配的<之后不是/；在>之前加上环视(?:<!/)，也是如此；剩下的内容仍然由[<^>]+匹配。这样结果清晰了很多，三个部分互不干扰，所以整个表达式就是<(?:!/)[<^>]+(?:<!/)>。

环视的另一价值在于，提取数据时杜绝错误的匹配。比如匹配邮政编码，直接的想法是找到6位数字构成的字符串，但仅仅用\d{6}提取，很可能在手机号码13812345678、电话号码28812506等其他数据中找到6位数字构成的字符串。如果在表达式首尾添加环视，改为(?:<!\\d)\\d{6}(?!\\d)，就可以保证准确匹配6位数字构成的字符串。一般来说，凡是提取“有长度特征的数据”，都需要用到环视。

还有些时候，可以在匹配的同时以环视施加限制，达到“双管齐下”的效果。举个例子，Java和.NET的正则表达式都提供了字符组运算的功能，比如匹配所有的辅音字母，最简单的思路是“从26个字母中减去5个元音字母”，在.NET中可以写作[a-z-[aeiou]]，在Java中可以写作[[a-z]&&[^aeiou]]。如果使用其他语言，则没有这种便利，只能写作[b-df-hj-np-tv-z]，不但烦琐，而且不便理解。但是使用环视则非常简单，写作(?:[aeiou])[a-z]，[a-z]真正匹配的是一个小写字母，但环视(?:[aeiou])同时要求这个字母不能由[aeiou]匹配，最终效果就是“从26个字母中减去5个辅音字母”。

### 4.4.2 环视与分组编号

环视结构也要用到括号，这种括号是否会影响到分组编号呢？

前面说过，分组的编号只与捕获型括号有关，而不受其他任何类型括号的影响。所以，环视结构虽然出现了括号，但并不影响分组。

#### 例 4-21 环视结构并不影响引用分组

```
#表达式为(?!=ab)cd
print re.search(r"(?!=ab)(cd)", "abcd").group(0)
cd

print re.search(r"(?!=ab)(cd)", "abcd").group(1)
cd
```

但是，括号有多种用途，比如表示多选结构。即便用于表示多选结构的括号，如果没有显式指定为非捕获型括号 `(?:...)`，也会视为捕获型括号，这时候结果就大不一样了。

#### 例 4-22 环视结构中出现了捕获型括号，会影响分组

```
#环视结构中出现了捕获型括号
print re.search(r"^(?=(ab|cd))", "abcd").group(0)

print re.search(r"^(?=(ab|cd))", "abcd").group(1)
ab

#环视结构指定使用非捕获型括号
print re.search(r"^(?=(?:ab|cd))", "abcd").group(1)
Traceback (most recent call last):
IndexError: no such group
```

这一点在实际使用中往往容易忽略，认为环视结构中的表达式匹配的文本并不能从匹配结果中得到，结果算错了真正需要的分组编号。

而且，环视结构中的捕获型括号一旦匹配完成，就不能回溯。在某些情况下，这可能带来非常奇怪的问题。比如表达式 `(\d+)\w+\1`，它可以匹配 `123a12`，其中 `\d+` 匹配 `12`，`\w+` 匹配 `3a`，`\1` 反向引用 `\d+` 匹配的 `12`。但是，将这个表达式改为 `(?<=(\d+)\w+\1)`，却无法在 `123a12` 中找到匹配，因为一开始 `\d+` 就会匹配 `123`，之后跳出环视结构，这时 `\d+` 保存的备用状态全部消失了，所以无法交还 `3` 给 `\w+` 匹配，导致整个表达式匹配失败。

不过，需要用到这种表达式的情况很罕见，而且环视中能出现的表达式往往会有限制，不见得能出现括号。

### 4.4.3 环视的支持程度

常用的语言大都支持环视，但语言不同，支持的程度也不同，下面详细讲解。

一般来说，所有语言都支持两种顺序环视，而且没有限制。也就是说，无论你使用肯定顺序环视，还是否定顺序环视，都可以在其中使用各种复杂的表达式。

逆序环视的情况则没有这么乐观，JavaScript 和 Ruby 1.8 中的正则表达式并不支持逆序环视，



其他语言虽然支持逆序环视，但对逆序环视中的表达式能匹配的文本长度有限制<sup>1</sup>：Python 只支持匹配固定长度文本的表达式，而 Java 和 PHP 只支持匹配有限长度文本的表达式，只有 .NET 没有任何限制。下面我们依次讲解这几种情况以及应急的解决办法。

### JavaScript、Ruby 1.8

JavaScript 和 Ruby 1.8 不支持逆序环视。如果确实需要用到逆序环视，而且只用到肯定逆序环视，不妨采用分组来替代，比如把表达式 `(?<=dog)s` 改写为 `dog(s)`。注意：我们给 `s` 添加了一个括号，如果使用 `(?<=dog)s`，需要提取整个表达式匹配的文本，但使用 `dog(s)`，只需要提取编号为 1 的分组捕获的文本即可。当然也可以更进一步，对 `dog` 使用非捕获型括号，写成 `(?:dog)s`，思路与 `dog(s)` 相同，却不需要改动原有的分组编号。

但是，这个办法只对肯定逆序环视有效，而不能用于否定逆序环视，因为逆序环视是“从右向左”判断的，其他结构都是“从左向右”判断的，如果逆序环视中表达式能匹配的字符串长度是不固定的，就很难确定“从左向右”判断时的起点。

即便否定逆序环视中表达式能匹配的字符串长度是固定的，也难以用其他结构解决。比如表达式 `(?!dog)s`，要求 `s` 左侧不能出现 `dog`，或许你会想用表达式 `(?!dog)(?:.{3})s`，但是如果字符串 `s` 或者 `is` 也无法匹配，因为这时候 `s` 左侧或者是空字符串，或者是只包含字母 `i` 的字符串，而表达式中 `s` 左侧的 `.{3}` 要求匹配三个字符。看来，表达式要改为 `(?!dog)(?:.{0,3})s`，但是这时候，它又可以匹配 `dogs` 中的 `ogs` 或者 `gs`……总的来说，没有妥善的解决办法。

### Python、Ruby 1.9

Python 规定在逆序环视中的表达式能匹配的文本长度必须是固定的。也就是说，`(?<=dog)` 是合法的，`(?<=(dog|cat))` 也是合法的，因为环视中的子表达式能匹配的文本都是固定的；而 `(?<=dogs?)` 和 `(?<=(dog|cats))` 则都不合法，因为环视中的子表达式能匹配的文本长度不确定。

要解决这个问题，可以用多选结构来改造表达式。比如，`dogs?` 等价于 `(dog|dogs)`，所以 `(?<=dogs?)` 可以改写为 `((?<=dog)|(?<=dogs))`，同样，`(?<=(dog|cats))` 可以改写为 `((?<=dog)|(?<=cats))`；但是这种办法的应用场景有限，如果逆序环视中的表达式比较复杂，用多选结构列出就非常麻烦；而且一旦表达式中出现了 `*` 和 `+` 之类的量词，就根本不可能用多选结构列出了。

<sup>1</sup> 这些限制之所以存在，因为逆序环视的机制与正则表达式匹配的机制完全不同，它从当前位置开始，由右向左“倒过来”查找可能的匹配。实际的匹配过程更像从右向左截取一段文本，再测试它能不能由表达式匹配。如果表达式能匹配的文本长度是固定的，处理的代价就很小，否则代价可能很大。

## PHP

PHP 对逆序环视中表达式的限制要宽松点，它能匹配文本的长度可以不必限定，但是长度必须是确定的数值：`(?<=dog|cats)`是没问题的，因为两个分支的长度都是固定的；而`(?<!dogs?|cats?)`则有问题，因为两个分支的长度都是不确定的。

要解决这个问题，可以将`(?<!dogs?|cats?)`改写为`(?<!dog|dogs|cat|cats)`。但是要注意，根据 PHP 中正则表达式的规定，长度固定但不相同的多选分支，只能出现在顶级的多选结构中。也就是说，`(?<=ab(c|de))`是不支持的，而`(?<=abc|abde)`则是支持的。

## Java

Java 对逆序环视中表达式的限制比 PHP 还要宽松，它能匹配的文本长度可以没有确定长度，但是必须有上限。所以`(?<!dogs?|cats?)`是没有问题的，甚至`(?<!ab(c|de))`也是没有问题的，但是`(?<!(dogs?){3,})`则不行，编译时会报告无法确定环视中表达式能匹配文本的最大长度。

## .NET

.NET 对逆序环视的支持是最为完备的，在.NET 的逆序环视中，表达式没有任何限制，这是.NET 的正则表达式最为人称道的一点。

总的来看，语言不同，对逆序环视的限制也不相同。这些限制之所以存在，是因为逆序环视的机制与正常的匹配机制完全不同，它从当前位置开始，由右向左“倒过来”查找可能的匹配。实际的操作过程更像从右向左截取一段文本，再判断它能不能由表达式匹配，这样的过程可能要重复尝试很多次。如果表达式能匹配的文本长度确定，处理的代价就很小，否则代价可能很大。所以比较好的做法是，即便支持，也应当尽量避免在逆序环视中使用复杂的表达式。

### 4.4.4 环视的组合

环视匹配的并不是字符，而是位置。在正则表达式匹配时，环视结构匹配成功，并不会更改“当前位置”，所以多个环视可以组合在一起，实现在同一个位置的多重判断。

最常见的组合是**环视中包含环视**（逻辑如图 4-4 所示），比如之前在匹配主机名时，我们限定主机名的长度不能超过 255 个字符，使用表达式`(?=[-a-zA-Z0-9.]{0,255})(?![-a-zA-Z0-9.])`。其中`(?![-a-zA-Z0-9.])`是包含在外层的环视中的，它要求在这个位置（也就是主机名字符串之后）不能再出现属于主机名字符串的字符，也就是保证之前的表达式匹配整个主机名字符串，而不是“可能的主机名字符串的一部分”；综合起来，`(?=[-a-zA-Z0-9.]{0,255})(?![-a-zA-Z0-9.])`保证的是“整个主机名字符串的长度在 255 个字符以内”。

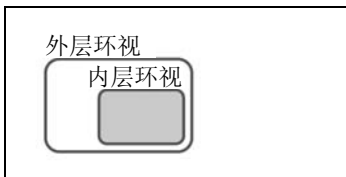


图 4-4 环视中包含环视

另一种组合是**并列多个环视**（逻辑如图 4-5 所示），它要求在当前位置，所有环视的判断都必须成功。比如要找到这样的位置：它之后是一个数字字符串，但不能是 999 开头的数字。这时候，就必须并列两个环视。表示数字字符串的表达式是 `\d+`，对应的环视结构是 `(?=\d+)`；表示“不是 999 开头”的表达式是 `(?!999)`。现在要做的是把两个环视并列起来，得到 `(?=\d+)(?!999)`。因为环视结构不会更改当前位置，所以先后顺序无所谓，无论是 `(?=\d+)(?!999)` 还是 `(?!999)(?=\d+)`，效果是相同的，都是要求下面两个条件同时满足：在当前位置，之后必须出现数字字符串；在当前位置，之后不能出现 999。最终的结果都是对两个环视做“与（and）”运算，也就是说，两个条件必须同时满足才算匹配成功，否则宣告当前位置匹配失败<sup>1</sup>。

```
#查找这样的起始位置：它之后是数字字符串，不过不能以 999 开头
#数字字符串
re.search(r"^(?=\d+)(?!999)", "1234") != None    # => True
#非数字字符串
re.search(r"^(?=\d+)(?!999)", "abcd") != None    # => False
#虽然是数字字符串，但是 999 开头
re.search(r"^(?=\d+)(?!999)", "9991234") != None # => False
```

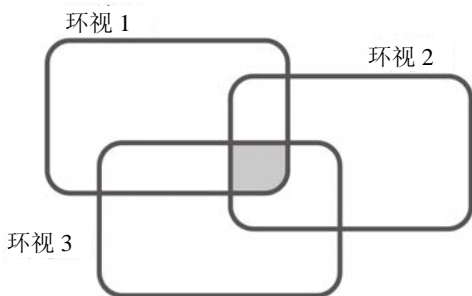


图 4-5 多个环视并列

最后一种组合是将若干个**环视作为多选分支排列在多选结构中**（逻辑如图 4-6 所示）。比如

<sup>1</sup> 如果仔细观察可能会发现，字符串 9999 的开始位置也不能匹配，这是因为 `(?!999)` 不容许出现的其实是字符串 999，而不是数值 999，如果要更准确地表示数值 999，应该使用 `(?!999(?:!\d))`。

要找到这样的位置：它之后要么不是数字字符，要么是一个数字字符和一个非数字字符（比如 **1a**）。“不是数字字符”对应的环视是 `(?!\d)`；而“一个数字字符和一个非数字字符”对应的环视是 `(?=\d\D)`，所以总的环视就是 `((?!\d)|(?=\d\D))`。虽然上一段也提过并列多个环视的组合，但在多选结构中列出多个环视结构的意义大不一样。使用多选结构时，列出的多个环视只要有一个成立，整个判断就成功；不使用多选结构时，所有列出的环视都必须成立，整个判断才成功，具体的例子可见例 4-24。

#### 例 4-23 环视作为多选分支

```
#查找这样的起始位置，之后要么不是数字字符，要么是一个数字字符和一个非数字字符
#不是数字字符
re.search(r"^((?!\d)|(?=\d\D))", "ab") != None    # => True
#是一个数字字符和一个非数字字符
re.search(r"^((?!\d)|(?=\d\D))", "a4") != None    # => True
#单个数字字符
re.search(r"^((?!\d)|(?=\d\D))", "4") != None     # => False
#连续数字字符
re.search(r"^((?!\d)|(?=\d\D))", "44") != None    # => False
```

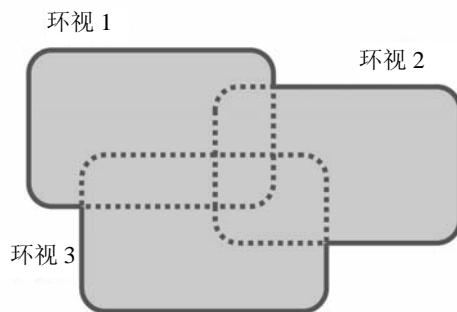


图 4-6 将环视结构排列在多选分支内

### 4.4.5 断言和反向引用之间的关系

断言不匹配任何字符，只匹配位置；而反向引用只引用之前的捕获分组匹配的文本，之前捕获分组中锚点表示的位置信息，在反向引用时并不会保留下来。

举例来说，如果表达式是 `(\bcat\b)\s+\1`，`\1` 所匹配的，就不只有单独出现的 **cat**，还包括单词内部的 **cat**（比如 **cate** 中的 **cat**），如果要验证单词 **cat** 是否在字符串中出现了两次，正确的做法是在反向引用两端也加上单词边界 `\b`，变成 `(\bcat\b).*?\b\1\b`，代码见例 4-25。

#### 例 4-24 反向引用时不会保留断言的判断

```
re.search(r"(\bcat\b).*?\1", "cat cate") != None    # => True
#两端添加\b之后，原有字符串无法匹配
```

```
re.search(r"(\bcat\b)\s+\b1\b", "cat cate") != None # => False
re.search(r"(\bcat\b)\s+\b1\b", "cat cat") != None # => True
```

有些正则表达式的资料中提到用反向引用匹配重复单词的例子，但是往往只使用 `(\b\w+\b)\s+\1` 来匹配“重复单词”，这其实是不对的，而应该使用 `(\b\w+\b)\s+\b\1\b`。反向引用时，之前捕获分组中的断言都会被忽略，这一点请务必注意。

## 第 5 章 匹配模式

前几章已经介绍了字符组、量词、括号、断言，本章讲解正则表达式的另一个常用功能：匹配模式。

所谓**匹配模式**（match mode）<sup>1</sup>，指的是匹配时使用的规则。设置特定的模式，可能会改变对正则表达式的识别，也可能会改变正则表达式中字符的匹配规定。常用的匹配模式一共有 4 种：**不区分大小写模式**、**单行模式**、**多行模式**、**注释模式**。

### 5.1 不区分大小写模式

通常，用户关心的只是文本的意义，而不是它的具体形式。比如单词 `the`，在句子中写作 `the`，在句子开头写作 `The`，还可能为了强调写作 `THE`；可是用户不管这些，只希望找到所有的 `the`。为达到这个目的，可以使用第 1 章介绍的字符组，写作 `[tT][hH][eE]`，这样做确实没错，但是如果单词长一些，写起来就很麻烦了，比如 `tomorrow` 就要写成 `[tT][oO][mM][oO][rR][rR][oO][wW]`。更重要的是，这样的表达式不够直观，读代码的人很难明白 `[tT][oO][mM][oO][rR][rR][oO][wW]` 要匹配的是 `tomorrow`。

为解决这种问题，正则表达式提供了**不区分大小写**<sup>2</sup>的匹配模式，指定此模式之后，在正则表达式中可以直接写 `the`，就可以匹配 `the`、`The`、`THE` 等各种大小写形式的 `the`，`tomorrow` 也可以匹配各种形式的 `tomorrow`，大大降低了理解的难度。

在看这个模式的应用实例之前，必须首先了解模式的指定方式。通常，有两种办法指定匹配模式：以**模式修饰符**指定，或者以**预定义的常量**作为特殊参数传入来指定。

<sup>1</sup> 此处“模式”对应的英文是 `mode`，而不是 `pattern`，许多语言的文档中把正则表达式本身叫做 `pattern`，请注意区别。

<sup>2</sup> `Case-Insensitive`，有许多人翻译为“不敏感”，这是不对的。按照词典的解释，“敏感/对某些方面感知特别敏锐”只是 `sensitive` 的意思之一，`sensitive` 还可以表示“有能力分辨/分别”，显然 `Case-Insensitive` 应当取这个意思，所以翻译为“不区分大小写”。

5.1.1 模式的指定方式

模式修饰符即模式对应的单个字符，使用时将其填入特定结构`(?modifier)`中（其中的`modifier`为模式修饰符），嵌在正则表达式的开头。比如不区分大小写的匹配模式对应的模式修饰符是`i`（`case Insensitive`），对`the`指定此模式，完整的正则表达式就是`(?i)the`。

这个表达式几乎可以原封不动地用在任何语言中（只有 JavaScript 是例外，JavaScript 不支持模式修饰符`(?modifier)`的记法，详见第 12 章），其意义都是相同的：对 `the` 采用不区分大小写的模式。Python 中的结果如例 5-1 所示。

例 5-1 用模式修饰符指定不区分大小写模式

```
re.search(r"(?i)the", "The") != None    # => True
re.search(r"(?i)the", "the") != None    # => True
re.search(r"(?i)the", "THE") != None    # => True
```

另一种指定模式的方式是使用预定义的常量作为参数，传入正则函数。在 Python 中不区分大小写的预定义常量是 `Re` 类的静态成员 `re.IGNORECASE`（一般来说，它都是某个类的静态成员），在 Java 中，属于 `Pattern` 类，在 .NET 中属于 `RegexOptions`，在 Ruby 中属于 `Regexp`。PHP 和 JavaScript 是例外，它们的做法（其实在 Ruby 中也可以这么做）是在正则表达式末尾的分隔符（`delimiter`）之后加上模式对应的字母（比如不区分大小写模式对应的字母是 `i`，则添加字母 `i`）。表 5-1 列出了常用语言中的写法。

表 5-1 常用语言中不区分大小写模式的预定义常量

语言	常量
.NET	<code>RegexOptions.IgnoreCase</code>
Java	<code>Pattern.CASE_INSENSITIVE</code>
JavaScript	<code>/regex/i</code>
PHP	<code>/regex/i</code>
Python	<code>re.I</code> <code>re.IGNORECASE</code>
Ruby	<code>Regexp::IGNORECASE</code> <code>/regex/i</code>

例 5-2 各语言中用预定义常量指定不区分大小写模式<sup>1</sup>**.NET**

```
Regex pattern = New Regex("the", RegexOptions.IgnoreCase);
pattern.IsMatch("ThE"); // => True
```

或者

```
Regex.IsMatch("the", "ThE", RegexOptions.IgnoreCase); // => True
```

**Java**

```
Pattern pattern = Pattern.compile("the", Pattern.CASE_INSENSITIVE);
pattern.matches("ThE").find(); // => True
```

**JavaScript**

```
pattern = /the/i;
pattern.test("ThE"); // => True
```

**PHP =>**

```
preg_match("/the/i", "ThE") != 0; // => True
```

**Python**

```
re.compile(r"the", re.I).search("ThE") != None # => True
```

或者

```
re.search(r"the", "ThE", re.I) != None # => True
```

**Ruby**

```
(Regexp.new("the", Regexp::IGNORECASE) =~ "ThE") != nil # => True
```

或者

```
(/the/i =~ "ThE") != nil # => True
```

比较两种指定形式，模式修饰符较为通用，因为在各种语言中写法基本相同，而预定义常量在不同语言中写法不同。不过，两种形式的效果是相同的：无论以哪种方式，只要指定了不区分大小写模式，正则表达式在匹配时，就不会区分同一个字母的大小写形式（`(?i)the` 和 `(?i)THE` 是完全等价的）。

---

<sup>1</sup> 在使用预定义常量指定匹配模式时，往往并不是在与匹配直接相关的函数中使用预定义常量，而是需要生成正则表达式对象，之后才能匹配。比如在 Java 的例子中，就是生成一个 `Pattern` 对象，然后调用该对象的对应方法（`matches()`），得到一个 `Matcher` 对象，再对 `Matcher` 对象调用 `find()` 方法，才能真正完成匹配。而之前我们看到的 Python 的例子，都是直接调用函数，将正则表达式和字符串传入（比如 `re.search(r"the", "THE")`）。Java 的这种办法，是正则表达式应用中的“面向对象式”处理；而之前看到 Python 中的方法，是正则表达式应用中的“函数式”处理，在上面的例子中，PHP 用到了此方法；而 .NET、Python、Ruby、JavaScript 可以同时使用两种处理方法。关于正则表达式的处理方式，详见第 101 页。



之前看到过匹配 HTML 中 tag 的例子，比如匹配超链接 tag 的正则表达式、匹配图片 tag 及网页标题 tag 的正则表达式（[图 46](#)），虽然 XHTML 规范推荐 tag 名都用小写字母，但类似<IMG>的 tag 也很可能出现，为同时兼容大写字母，可以使用不区分大小写模式，这样比用字符组详细列出字母的大小写形式简单很多，如表 5-2 所示。

表 5-2 采用不区分大小写模式，简化匹配 tag 的表达式

任务	表达式
提取超链接	(?i)<a\s+href\s*=\s*["']?([^"'\s]+)["']?>([^\<]+)</a>
提取标题	(?i)<head>([^\>]+)</head>
提取图片	(?i)<img\s[^\>]*?src=["']?([^"'\s]+)["']?([^\>]*)>

5.2 单行模式

元字符点号 `.` 几乎能匹配任何字符，唯有换行符 `\n` 是例外。但是，有时候确实需要匹配“任何字符”，比如在处理 HTML 源代码时，经常会遇到跨越多行的脚本代码。

```
<script type="text/javascript">
...code...
...code...
</script>
```

正则文档里一般都会说明“点号不能匹配换行符”，不过许多人并不习惯仔细阅读文档，所以认为点号 `.` 能匹配任何字符，当然也就包括换行符，所以直接的想法是用 `<script\s.*?</script>` 来匹配。

如果你仔细读过前面的章节就会知道，因为这段 JavaScript 代码中出现了换行符，所以 `. * ?` 的匹配最多只能进行到第一行末尾。之前提到过，可以用 `[\s\S]` 之类的字符组匹配“任意字符”，所以正则表达式 `<script\s[\s\S]*?</script>` 能解决问题。

不过对大多数人来说，点号更加自然，也更简洁，所以正则表达式提供了**单行模式**。在这种模式下，所有文本似乎只在一行里，换行符是这一行中的“普通字符”，所以可以由点号 `.` 匹配。

单行模式对应的模式修饰符是 **s** (**S**ingle line)，所以如果用模式修饰符，可以在表达式的开头用 `(?s)` 指定，因此上面的表达式也可以改写为 `(?s)<script\s.*?</script>`。使用预定义常量的写法见表 5-3。

表 5-3 常用语言中单行模式的预定义常量

语言	常量
.NET	RegexOptions.Singleline
Java	Pattern.DOTALL

(续表)

语言	常量
JavaScript	不支持此模式
PHP	<code>/regex/s</code>
Python	<code>re.S</code> <code>re.DOTALL</code>
Ruby	<code>Regexp::MULTILINE</code> <code>/regex/m</code>

你可能注意到了，单行模式在不同语言中的称呼很不一样，甚至在同一门语言中，也可能有不同的记法（比如在 Python 中）。比如在 Java 和 Python 中叫做 DOTALL（也就是点号通配），这个名字确实更高明（因为常用的模式中还包含“多行模式”，它和“单行模式”没什么联系，但是这两个名字确实很迷惑人）。不过，“单行模式”成了约定俗成的称呼（通用的模式修饰符是 `s`，它难以联系上 DOTALL），所以本书还是沿用“单行模式”的说法。

比较奇怪的是 Ruby 的预定义常量 `Multiline`，它的意思是“多行”，而且它使用的模式修饰符也是 `m`。这确实让人费解，或许本意是“使用点号 `.` 的正则表达式可以跨越多行”？不管怎样，请一定记住，Ruby 中的“多行模式”实际上是常说的“单行模式”。

如果不想使用“点号+单行模式”的组合（比如 JavaScript 完全不支持这种模式，想用也没办法），也可以使用 `[\s\S]` 之类的字符组，它的确可以匹配任何字符，只是许多人并不习惯这样的写法，“点号+单行模式”的组合看起来更顺眼一些。

### 5.3 多行模式

“多行模式”听起来是与“单行模式”对应的，其实这两个模式没有任何联系。

**单行模式影响的是点号的匹配规则：**在默认模式下，点号 `.` 可以匹配除换行符之外的任何字符，在单行模式下，点号 `.` 可以匹配包括换行符在内的任何字符；**多行模式影响的是 `^` 和 `$` 的匹配规则：**在默认模式下，`^` 和 `$` 匹配的是整个字符串的起始位置和结束位置，但在多行模式下，它们也能匹配字符串内部某一行文本的起始位置和结束位置。

假设，需要找到下面文本中所有数字字符开头的行。

```
1 line
No digit
2 line
```

解决这个问题，需要定位到每行的起始位置，尝试匹配一个数字字符，如果成功，则匹配之后的整行文本。多行模式的模式修饰符是 `m` (**Multiline**)，所以在表达式的开头用 `(?m)` 指定多行模式，这样 `^` 可以定位到字符串内部每一行的起始位置；匹配数字字符的表达式是 `\d`，因为没有

指定单行模式，点号 `.` 不能匹配换行符，`.*` 可以匹配 “之后的整行文本”，整个表达式就是 `(?m)^\d.*`。

例 5-3 用模式修饰符指定多行模式

```
multilineString = "1 line\nNot digit\n2 line"
lineBeginWithDigitRegex = r"(?m)^\d.*"
for line in re.findall(lineBeginWithDigitRegex, multilineString):
    print line

1 line
2 line
```

还可以利用多行模式下的 `$`，给每一行的末尾添加句号。

例 5-4 在多行模式下，给行末添加句号

```
multilineString = "1 line\nNot digit\n2 line"
lineEndRegex = r"(?m)$"
print re.sub(lineEndRegex, ".", multilineString)

1 line.
Not digit.
2 line.
```

上面的表达式几乎是任何语言中都“通用”的，唯有 JavaScript 是例外，因为它不支持用模式修饰符指定模式，但是可以使用预定义常量来指定多行模式。在表 5-4 中列出常用语言中预定义常量的写法。

表 5-4 常用语言中多行模式的预定义常量

语言	常量
.NET	RegexOptions.Multiline
Java	Pattern.MULTILINE
JavaScript	/regex/m
PHP	/regex/m
Python	re.M re.MULTILINE
Ruby	默认即为多行模式

在各种语言中，多行模式对应预定义常量的写法比较统一，都是 `multiline`。值得一提的是 Ruby，它默认就采用多行模式，`^` 和 `$` 在任何情况下都能匹配文本内部的行起始/结束位置。如果要在 Ruby 中“摆脱”多行模式，只能用 `\A` 替代 `^`，用 `\Z` 替代 `$`。

## 5.4 注释模式

有时，用到的正则表达式可能非常复杂，不但难以编写和阅读，也难以维护。如果正则表达式也像程序源代码一样，可以添加注释，阅读和维护就容易多了。

为解决这个问题，许多语言支持使用 `(?#comment)` 的记法添加注释，`comment` 就是注释的内容。所以，上面的表达式 `^\d.*?$` 就可以写成这样：

```
^(?#start of the line)\d(?#digit).*(?#rest of the line)
```

.NET、Python、Ruby、PHP 都支持这种记法，Java 和 JavaScript 则不支持。不过，还有一种注释的写法是各种语言都支持的，就是使用注释模式，此时，正则表达式对应的字符串可以跨越很多行<sup>1</sup>。注释模式的代码见例 5-5。

例 5-5 注释模式

```
multilineString = "1 line \nNot digit\n2 line"
lineBeginWithDigitRegex = r"""
(?mx)      #enable multiline and extended mode
^          #start of the line
\d         #digit
.*         #rest of the line
$          #end of the line
"""

print re.findall(lineBeginWithDigitRegex, multilineString)
['1 line', '2 line']
```

在注释模式下，正则表达式内部的空白字符都被忽略（一般来说，主要是 ASCII 编码中的空白字符，Unicode 编码中的空白字符情况不定），注释则以 `#comment` 的形式添加在正则表达式内部，每一条注释从 `#` 开始，到行末结束。在许多文档中，都是用这种模式来解释复杂的表达式，并且会使用缩进表示层级结构，这样更加方便阅读和维护。比如匹配日期的正则表达式 `((?x)(\d{4})-(\d{2})-(\d{2}))`，在注释模式下可以就这样展开。

例 5-6 在注释模式下展开复杂正则表达式

```
dateRegex = r"""
(?x)          #enable multiline and extended mode
(             #start of whole regex
    (\d{4})    #year
    -         #dash
    (\d{2})    #month
)
```

<sup>1</sup> 这里为了方便展示，使用了 Python 中的三引号字符串（Triple-Quoted String），这是 Python 中的特殊写法，以三个引号为首尾标志，字符串的值完全就是期间的文本“看上去的样子”，所以这个字符串横跨了多行。但是三引号字符串中也会进行字符串转义，所以在之前添加了 `r`，表示原生字符串。

```
-          #dash
(\d{2})    #day
)          #end of whoe regex
""""
```

注释模式对应的模式修饰符是 **x** (**eXtended mode**，扩展模式，但更常见的写法是 **free-spacing mode**，**宽松格式模式**)。表 5-5 介绍了各语言中注释模式的预定义常量。

表 5-5 各语言中注释模式的预定义常量

语言	常量
.NET	RegexOptions.IgnorePatternWhitespace
Java	Pattern.COMMENTS
JavaScript	/regex/x
PHP	/regex/x
Python	re.X re.VERBOSE
Ruby	Regexp::EXTENDED /regex/x

你可能注意到了，例 5-6 同时指定了两种模式：多行模式和注释模式。注释模式的 **x** 与多行模式的模式修饰符 **m**，合写作 **(?mx)**。如果需要同时使用多种模式，只要在 **(?modifier)** 中将模式修饰符排列起来就可以了。

如果希望同时使用多行模式和注释模式，使用预定义常量该怎么做？答案是，使用位运算符 **|**。通常来说，匹配模式对应的预定义常量都是 **int** 类型，所以多个值进行按位与的结果，并不会彼此干扰。比如在 **Java** 中，对应的写法就是 **Pattern.COMMENTS | Pattern.MULTILINE**，在 **.NET**、**Ruby** 和 **Python** 中也可以这样。

如果是 **PHP** 和 **JavaScript**，则在结束的分隔符之后直接并列模式对应的修饰符，比如 **/regex/mx**，这种写法在 **Ruby** 中也行得通。

## 5.5 补充

### 5.5.1 更多的模式

上面提到的四种常用模式是各语言中通用的，但是匹配模式并不只有这 4 种，不同的语言提供了不同的模式，它们一般都在预定义常量中，而且不一定有对应的模式修饰符。

比如 Java 的 `Pattern` 还包含其他模式（仅列举两种模式，更多的请参看文档）。

`Pattern.UNIX_LINES`: 在此模式下，`^`、`.`、`$` 识别的“行终止符”只有 `\n`，而不能是其他字符（比如 `\r\n`），它对应的模式修饰符是 `d`；

`Pattern.CANON_EQ`: 在此模式下，字符的“相等”规则更加灵活，Unicode 字符 `a\u030A` 与 `\u00E5` 也是“相等”的（都是 `å`，也就是拉丁字母 `a` 加上分音符，只是前者用两个字符组合，后者用单个字符），此模式可能对性能有很大影响，而且没有对应的模式修饰符。

Python 的 `re` 也包含了其他模式（仅列举两种模式，更多的请参阅文档）。

`re.U` 或 `re.UNICODE`: 在此模式下，`\w`、`\d`、`\s` 等字符组简记法的匹配规则会发生改变，比如 `\w` 能匹配 Unicode 中的“单词字符”，包括中文字符，`\d` 也能匹配 `1`、`2` 之类的全角数字字符，它有对应的模式修饰符 `u`。

`re.A` 或 `re.ASCII`: 因为在 Python 3 以上的版本中，正则表达式默认采用 Unicode 匹配规则，如果希望让 `\d`、`\w` 等字符组简记法恢复到 ASCII 匹配规则，可以使用此模式，它有对应的模式修饰符 `a`。

在 .NET 和 PHP 中，也有其他模式可用，不过一般来说，本章介绍的 4 种模式最为常用，在其他模式中，只有涉及 Unicode 或者 ASCII 的模式使用较多，因为它们影响了字符组简记法 `\d`、`\s`、`\w` 的匹配规则。关于每种语言可用模式的具体信息，请参考该语言对应的章节或文档。

5.5.2 修饰符的作用范围

常见的模式修饰符是 `(?modifier)` 的形式，它表示“从现在开始使用某个模式”。通常做法是将它写在正则表达式的最开头，表示“整个正则表达式都指定此模式”；如果它出现在正则表达式当中，则表示 **此模式从这里开始生效**<sup>1</sup>；如果模式修饰符出现在某个括号内——比如 `((?modifier)...)` ——那么它的作用范围只限于括号内部，此模式也可以记为 `(?modifier:...)` 。

正则表达式	能匹配的文本
<code>t(?i)he</code>	<code>tHE the thE tHe</code>
<code>th(?i)e</code>	<code>thE the</code>
<code>t((?i)h)e</code>	<code>tHe the</code>
<code>t(?i:h)e</code>	<code>tHe the</code>

模式修饰符对正则表达式的操控性更强，因为预定义常量指定的匹配模式是对整个表达式生效的。

<sup>1</sup> Python 的情况有所不同，模式修饰符 `(?modifier)` 只要出现，无论在什么位置，都对整个正则表达式生效。

### 5.5.3 失效修饰符

`(?modifier)` 指定了匹配模式开始作用的范围，在正则表达式中，另有一类失效修饰符，它用来“结束”某种模式的作用范围，其形式是 `(?-modifier)`，类似 `(?modifier)`，只是问号 `?` 之后多了一个减号 `-`，表示“取消模式”，也就是 **某个模式生效到此为止**。

假设，一段文本中包含了许多单词，其中以 `bar` 结尾的有 `foobar`、`zeebar` 等，而且大小写不定。

```
FooBAR feeBAR zeeBAR ZEEBAR foobar zeeBAR
```

现在希望找出所有这样的 `foobar` 和 `zeebar`：不论 `foo` 和 `zee` 的大小写如何，只要结尾是大写的 `BAR` 就可以。`FooBAR`、`ZEEBAR`、`zeeBAR` 都符合要求，`FooBAR`、`zooBAR`、`feeBAR` 等则要排除。

你可能会觉得这很简单，要找到的单词，前面部分可能是 `zee` 或是 `foo`，不区分大小写则用字符组把大小写形式都列出来即可，后面必定是 `BAR`；所以正则表达式是 `[fFzZ][oOeE][oOeE]BAR`。可是，这个表达式行不通，因为它可以匹配 `feeBAR`，而这并不是我们需要的。

真正要做的，其实是准确划定不区分大小写模式的作用范围即可：匹配前面部分的表达式是 `(foo|zee)`，同时必须使用不区分大小写模式；匹配后面部分的 `BAR` 则必须停止使用不区分大小写模式。

要满足这个要求，可以使用上一节介绍的 `((?i)foo|zee)BAR`，也可以使用失效修饰符，将表达式写作 `(?i)(foo|zee)(?-i)BAR`，代码见例 5-7，因为 Python 不支持这种写法，所以使用了 Ruby（Ruby 的正则使用细节请参考第 15 章）。

#### 例 5-7 不区分大小写模式与反向引用

```
str = " FooBAR feeBAR zeeBAR ZEEBAR foobar zeeBAR"
regex = /(?(i)(foo|zee)(?-i)BAR/
puts str.scan(regex)
["feeBAR", "ZEEBBAR", "zeeBBAR"]
```

一般来说，只要语言或工具支持模式修饰符 `(?modifier)`，都可以支持失效修饰符 `(?-modifier)`，两者配合使用，可以更精确地设定模式的作用范围。作为例外，Python 和 JavaScript 并不支持 `(?-modifier)` 的写法。

### 5.5.4 模式与反向引用

反向引用前面介绍过，它引用之前的表达式匹配的文本。如果之前的表达式使用了某种模式，引用时是否会继承这种模式呢？下面看看实际情况，因为在 Python 中无法限定模式的作用范围，现在以 Java 为例（具体语法请参考第 11 章），代码见例 5-8。

### 例 5-8 不区分大小写模式与反向引用

```
//不区分大小写模式
"abcABC".matches("(?i)abc\\1");    // => False
"abcABC".matches("(?i)(abc)\\1");  // => True
```

在表达式 `((?i)abc)\1` 中, `\1` 引用的是 `(?i)abc` 能匹配的文本(在这个例子里就是 `abc`), 因为 `\1` 不在区分大小写模式的作用范围内, 所以无法匹配 `ABC`; 而在表达式 `(?i)(abc)\1` 中, `\1` 处在区分大小写模式的作用范围内, 所以可以成功匹配 `ABC`。

### 例 5-9 单行模式与反向引用

```
//单行模式
"a\na\n".matches("(?s)a.\\1");    // => True
"a\na\n".matches("(?s)(a.\\1");  // => True
```

在例 5-9 的两个表达式中, `\1` 引用的是之前 `a.` 匹配的文本, 在单行模式下 `.` 可以匹配换行符; 所以 `a.` 匹配的是字符 `a` 和换行符 `\n`, 因此 `\1` 也可以匹配字符 `a` 和换行符 `\n`。

### 例 5-10 多行模式与反向引用

```
//多行模式
"aba".matches("(?m)^a.\\1");    // => True
"aba".matches("(?m)(^a.\\1");  // => True
```

在例 5-10 的两个表达式中, `\1` 引用的是之前 `^a` 匹配的文本, `^a` 只能匹配行开头的 `a`; 但是无论 `\1` 是否处在多行模式的作用范围内, 它能匹配的不只有行开头的 `a`, 而是任何位置的字符 `a`, 这一点值得注意。

## 5.5.5 冲突策略

使用 `(?-modifier)` 可以终止某个模式的作用范围; 但是, 模式也可以通过预定义常量来指定, 它是对整个表达式生效的。这时候, 就可能产生冲突。

如果在正则表达式中使用 `(?-i)` 取消不区分大小写的匹配模式, 又使用了预定义常量 (比如 Java 中的 `Pattern.CASE_INSENSITIVE`) 指定整个表达式采用不区分大小写的模式, 这时候情况会是怎样呢? 从例 5-11 可以看到 (因为 Python 不支持失效修饰符, 本节的例子均使用 Java, Java 的正则使用细节参见第 11 章)。

### 例 5-11 冲突策略

```
Pattern.compile("(?-i)a",
    Pattern.CASE_INSENSITIVE).matches("A").find(); // => False
```



可以看到，**模式修饰符具有更高的优先级**。严格地说，如果用预定义常量指定了整个正则表达式采用某种模式，同时正则表达式内部使用了关闭该模式的失效修饰符，则失效修饰符之后的部分，都不受预定义常量所指定模式的影响（失效修饰符之前的部分则不受影响），例 5-12 给出具体的例子。

例 5-12 更详细的例子

```
Pattern regex = Pattern.compile("a(?-i)a", Pattern.CASE_INSENSITIVE);
regex.matches("aa").find(); // => True
regex.matches("Aa").find(); // => True
regex.matches("aA").find(); // => False
```

### 5.5.6 哪种方式更好

模式修饰符和预定义常量都可以指定匹配模式，哪种方式更好？这个问题没有标准答案。

用模式修饰符指定的好处之一是简洁，因为通用的模式修饰符就是 `i x s m`，虽然不那么直观，但习惯之后并不难理解；而且，这些记法在各语言中都是通用的；再者，模式修饰符还可以用 `(?i)`、`(?-i)` 来精确控制模式的作用范围，这是预定义常量做不到的；另外，许多正则处理函数（详见第 6 章）只接受字符串形式的正则表达式，此时只能以模式修饰符来标识模式。

使用预定义常量的好处在于，它很形象——“单行模式”、“多行模式”的说法确实很让人困惑，尤其他们竟然可以同时使用，互不干扰，而看字面意思分明是互相矛盾的；而且，一些模式并没有对应的修饰符，比如 Java 中的 `Pattern.LITERAL`（消除所有元字符的特殊含义）、.NET 中的 `RegexOptions.ECMAScript`（字符组简记法采用 ASCII 匹配规则）。

我的建议是，熟练掌握常用的模式，在能够使用模式修饰符的地方，尽量使用修饰符。毕竟，无论使用哪种编程语言，`(?i)` 是熟悉正则表达式的人都能看懂的，但是 Java 中的 `Pattern.DOTALL` 等价于 .NET 中的 `RegexOptions.Singleline`、`Pattern.CASE_INSENSITIVE` 和 `RegexOptions.IgnoreCase`，理解起来总是要多点周折。

# 第 6 章 其他

前 5 章已经介绍了正则表达式常用的功能和结构，但在实际应用中，经常可能遇到一些繁杂具体的问题，而常见的文档中往往不会涉及这类问题。所以有必要用一章的篇幅讲解实用中经常遇到的问题，澄清若干容易混淆的概念。

## 6.1 转义

正则表达式中的转义是个麻烦的问题，经常会搞得人头疼。虽然之前各章都提到了转义，仍然有必要专门介绍。

### 6.1.1 字符串转义与正则转义

理解转义的基础是，明白字符串与正则表达式的关系。通常说的 *string*（字符串）中，*string* 称为**字符串文字**（String Literal）<sup>1</sup>，它是某个字符串的值在源代码中的表现形式。比如字符串文字 `\n`，它包含 `\` 和 `n` 两个字符，意义（或者说它的值）是一个换行符（为方便观察，表示为 `NL`）。在生成字符串时，应当进行“字符串转义”，才能准确识别字符串文字中 `\n` 的意思，如表 6-1 所示。

表 6-1 常见的字符串转义

字符串文字	字符串	说明
<code>\n</code>	<code>NL</code>	换行符
<code>\t</code>	<code>Tab</code>	制表符
<code>\\</code>	<code>\</code>	反斜线字符

另一方面，在源代码中看到的“正则表达式”*regex*，其中的 *regex* 称为**正则表达式文字**（Regular Expression Literal，以下简称正则文字），是正则表达式的表现形式。比如正则表达式 `\d`，其正则文字包含 `\` 和 `d` 两个字符，它的意义（或者说值）是匹配数字字符的字符组简记法。在生成正则表达式时应当进行“正则转义”，才能将正则文字中的 `\d` 识别为字符组简记法。

<sup>1</sup> 请参考 [http://en.wikipedia.org/wiki/String\\_literal](http://en.wikipedia.org/wiki/String_literal)。

这两点都不难理解，但从字符串转到正则表达式的过程就比较复杂。因为不少正则表达式都是以字符串形式提供的，必须经过从字符串文字到正则表达式的转换。根据上面的介绍，字符串文字首先必须经过“字符串转义”，才能得到真正的字符串；接下来，这个字符串作为正则文字，经过“正则转义”，最终得到正则表达式。图 6-1 说明了从字符串文字到正则表达式的转义过程，表 6-2 则列出了若干常用字符串的转义。

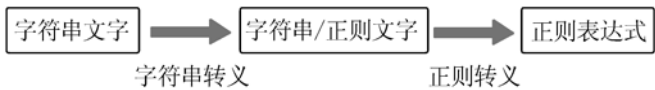


图 6-1 字符串转义与正则转义

表 6-2 以常见的几个字符为例，详细说明了这种转换过程。

表 6-2 从字符串文字到正则表达式的转换

字符串文字	字符串/正则文字	正则表达式	说明
\\n	\\n	<span>NL</span>	换行符
\\t	\\t	<span>Tab</span>	制表符
\\\\	\\	\\	反斜线字符

换行符和制表符的情况比较容易理解，它们的正则文字是\\n 和\\t，也就是字符串的值。但是从字符串文字必须经过字符串转义才得到作为正则文字的\\n 和\\t，根据字符串转义的规则，反斜线字符\\必须写成\\，所以字符串文字必须写成\\n 和\\t。

如果要表示正则表达式中的\\，必须使用正则文字\\，这样在正则转义时才能正确识别。同时，正则文字中的每个\\都必须用字符串文字\\表示，所以正则表达式\\对应的字符串文字就是\\\\。网络上常有人问，某个正则表达式中为什么需要用到 4 个反斜线字符，原因就在这里。

例 6-1 \ 的转义

```
#\\\\经过字符串转义和正则转义，得到的正则表达式就是一个\\
#\\经过字符串转义，得到的字符串就是一个\\
re.search("\\\\\\\\", "\\") != None      #  =>  True
```

有时候情况更复杂：正则表达式中的换行符或者制表符，在字符串文字中必须写成\\n 或\\t；但是，用\\n 或\\t 也没有问题。原因在于，在处理字符串转义时，它们已经被解释为换行符或制表符，所以传递给正则表达式的字符串中就包含了换行符或者制表符，具体情况如表 6-3 所示。

表 6-3 从字符串文字到正则表达式的转换

字符串文字	字符串/正则文字	正则表达式	说明
<code>\\n</code>	<code>\n</code>	<code>\\N</code>	换行符
<code>\n</code>	<code>\\N</code>	<code>\\N</code>	换行符
<code>\\t</code>	<code>\t</code>	<code>\\Tab</code>	制表符
<code>\t</code>	<code>\\Tab</code>	<code>\\Tab</code>	制表符

从例 6-2 可以看到，字符串文字中的`\\n`和`\n`表示的是同一个字符，只是前者先做字符串转义，再做正则转义，而后者只做了字符串转义。看来，似乎还是`\n`更符合直觉一些——如果你确实能够理解，为什么字符串文字中`\\n`和`\n`的结果竟然一样，使用`\n`确实可行。但是，正则表达式中`\`的字符串文字还是必须写成`\\`。

例 6-2 `\n`和`\t`的转义

```
re.search("\\n", "\n") != None # => True
re.search("\n", "\n") != None # => True
re.search("\\t", "\t") != None # => True
re.search("\t", "\t") != None # => True
```

要特别注意的是`\b`。在一般字符串中，`\b`是预定义的转义序列，表示退格符（backspace，表示为`BS`）；但是在正则表达式中，它表示单词边界（记为`\B`）。如果在字符串文字中写`\b`，字符串转义为退格符，作为正则表达式交给正则表达式，正则表达式真正得到的就是退格符`BS`，而不是单词边界`\B`，如表 6-4 和例 6-3 所示。所以，如果用到了单词边界，在字符串文字中一定要写成`\\b`。所以，最保险的办法是：正则表达式中的每一个`\`，在字符串文字中都要写成`\\`。本书之前章节中的代码，一直都遵循的正是这样的规则。

表 6-4 `\b`的转义

字符串文字	字符串/正则文字	正则表达式	说明
<code>\b</code>	<code>\\BS</code>	<code>\\BS</code>	退格符
<code>\\b</code>	<code>\b</code>	<code>\\B</code>	单词边界

例 6-3 `\b`的转义

```
re.search("\\ba\b", "a") != None # => False
re.search("\\\\ba\\b", "a") != None # => True
```

在 Python 中也可以这样。

例 6-4 Python 中`\d`的转义

```
re.search("\d", "1") != None # => True
re.search("\(", "(") != None # => True
```

看起来足够奇怪，因为在字符串中\nd 和\(\都不是合法的转义序列，所以如果在 Java 中这样写，是会报编译错误的。为什么 Python 中可以呢？这是因为 Python 对字符串有特殊的规定：如果遇到无法识别的转义序列，则将它原封不动地保存下来，具体的处理见表 6-5。

表 6-5 \nd 的转义

语言	字符串文字	字符串/正则文字	正则表达式	说明
Python	\\d	\d	<span style="border: 1px solid black;">\d</span>	字符组简记法
	\d	无法识别，原样保存	<span style="border: 1px solid black;">\d</span>	字符组简记法
Java	\\d	\d	<span style="border: 1px solid black;">\d</span>	字符组简记法
	\d	无法识别，报错	×	×

PHP 中也有同样的规定，所以在 PHP 的字符串文字中也可以直接写\d、\w、\s 等，效果和\\d、\\w、\\s 是一样的。

这样看来，使用字符串形式的正则表达式，转义的处理确实比较复杂。最好是能省去这些麻烦——正则表达式是怎样的，正则文字中就怎样写。要做到这一点，有两种办法：第一，使用原生字符串（Raw String），也就是完全忽略字符串转义的特殊字符串，在 Python、Ruby、.NET 中可以这么做。

```
Python
#在字符串之前用 r 标注原生字符串
pattern = re.compile(r"a\.\nb\b")

Ruby
#Ruby 中用单引号标注原生字符串
pattern = new Regexp('a\.\nb\b')

C#
//C#中用@标注
Regex regex = new Regex(@"a\.\nb\b");

VB.NET
//VB.NET 中的字符串与其他语言的不一样，只有双引号字符"需要转义
Dim regex as new Regex("a\.\nb\b");
```

第二，直接使用正则文字，在 Ruby 和 JavaScript 中可以这么做。下面以正则表达式 `a\.\nb\b` 为例。这种方法一般是直接在正则表达式两端添加分隔符`/`。此时，字符串转义会被略过，不过另外还有一个字符需要转义，就是分隔符`/`。

```
Ruby
pattern = Regexp.new(/a\.\nb\b/)
#下面这个表达式只包含一个字符，就是/
```

```
pattern = Regexp.new(/\\/))

JavaScript
var pattern = new RegExp(/a\.\\nb\\b/);
//下面这个表达式只包含一个字符，就是/
var pattern = new RegExp(/\\/));
```

总的来看，转义的情况确实很复杂。我推荐的做法是：

第一，如果可以使用原生字符串或者正则文字直接表示正则表达式，应当尽量使用这种做法，因为它简单直观，方便理解<sup>1</sup>；而且，本书从此开始，如果能使用原生字符串表示正则表达式，则使用原生字符串，在 Python 示例代码中使用 `r"regex"` 的形式。

第二，如果确定必须使用字符串文字，请尽量坚持这条原则：正则表达式中的每个反斜线`\`，在字符串文字中都必须写成`\\`，只有`\n`、`\t` 中的反斜线例外（但是`\\n` 和`\\t` 也不难理解）；我**不推荐**使用语言自身提供的“字符串转义时将无法识别的转义序列保留下来”的规定——比如要在正则表达式中使用字符组简记法`\d`，正则文字中要写为`\\d`，而不是`\d`——这样有利于使用其他语言的程序员阅读。

6.1.2 元字符的转义

了解了字符串与表达式的关系，再来看正则表达式中元字符的转义。元字符是有特殊含义的字符，如果要匹配“元字符”自身则必须转义，也就是在元字符之前添加反斜线`\`。比如元字符点号`.`，可以匹配除换行符外的任何字符，如果要准确匹配字符串中的点号`.`，正则表达式中就必须写`\.`。

也有些时候，匹配元字符自身并非一定要转义，表 6-6 总结列出了各种结构的转义。

表 6-6 常用结构的转义

结构	记法	转义	说明
字符组	[ ]	\[ ]	只对开方括号转义
	.	\.	
	-	\-	[a\~b]等价于[-ab]
量词	* + ?	\* \+ \?	
	*? +? ??	\*?\? \+?\? \?\?	
	{m,n}	\{m,n}	只对开花括号转义
括号	(...)	\(...\)	开、闭括号都需要转义

<sup>1</sup> 考虑到不是所有语言都支持原生字符串，所以本书的例子中，一般还是通过字符串来生成正则表达式，以利于读者理解，虽然有时候这样做确实很麻烦。

( 续表 )

结构	记法	转义	说明
多选结构	<code> </code>	<code>\ </code>	竖线需要转义
括号和多选结构	<code>( ...   ... )</code>	<code>\( ...   ... )</code>	开、闭括号和竖线都要转义
断言	<code>^ \$</code>	<code>\^ \\$</code>	
替换引用	<code>\$num</code>	<code>\\$或\$\$</code>	在替换的 <i>replacement</i> 字符串中转义

从表中可以看到，对称出现的元字符在转义时并不是“对称”的，比如与开方括号`[`对应的闭方括号`]`，与开花括号`{`对应的闭花括号`}`，这两个字符是否是元字符，取决于之前是否出现了开方括号或开花括号。如果出现了，则作为元字符出现；否则，作为普通字符出现。表 6-7 以方括号为例，说明了这一点。

表 6-7 字符组只要求转义开方括号

正则表达式	解释	能匹配的字符串
<code>[ab]</code>	字符组 <code>[ab]</code>	a b
<code>\[ab]</code>	字符 <code>[</code> 、字符 <code>a</code> 、字符 <code>b</code> 、字符 <code>]</code>	<code>[ab]</code>
<code>ab]</code>	字符 <code>a</code> 、字符 <code>b</code> 、字符 <code>]</code>	<code>ab]</code>

在字符组内部的闭方括号`]`在任何情况下都要转义，否则类似`[ ] ]`的正则表达式会出现二义性，造成识别错误，所以能匹配字符 `a`、字符 `b`、字符`]`的字符组，应当写为`[ab\ ]]`，而不是`[ab]]`。同样道理，括号内部的任何闭括号`)`都要转义，比如包含 `ab` 和 `b)` 的多选结构的正则表达式就应当写为`(ab|c\ ) )`，而不是`(ab|c) )`。

另一点容易忽略的是，在进行正则表达式替换时，*replacement* 字符串中也可能出现转义。比如在 Java 中，*replacement* 里通过`$num` 引用对应的捕获分组。如果想在替换时使用一个单独的`$`符号，而不是引用分组（比如生成价格字符串），则会报错，必须做转义才可以解决问题，在 Java 中使用`\$`，代码见例 6-5。

例 6-5 Java 中\$的转义

```
System.out.println("2010-12-20".replaceAll("(\\d{4})-(\\d{2})-(\\d{2})", "$2/$3/$1"));
12/20/2010

System.out.println("the price is 12.99".replaceAll("\\d+\\.\\d{0,2}", "$$0"));
Exception in thread "main" java.util.regex.PatternSyntaxException: Illegal
group reference
```

```
//转义需要出现的$
System.out.println("the price is 12.99".replaceAll("\\d+\\.\\d{0,2}", "\\$$$0"));
the price is $12.99
```

.NET 中的 *replacement* 中同样使用 *\$num* 引用对应的捕获分组，但是 .NET 中 *\$* 的转义并不是 *\\$*，而是 *\$\$*，代码见例 6-6。

#### 例 6-6 .NET 中 \$ 的转义

```
Console.Write(Regex.Replace("the price is 12.99", "\\d+\\.\\d{0,2}", "$$$0"));
the price is $12.99
```

### 6.1.3 彻底消除元字符的特殊含义

有些时候需要消除所有元字符的特殊含义，全部作为普通字符。这种情况经常发生在处理用户输入的场所：用户输入了某个字符串，需要根据这个字符串进行正则表达式查找。

假设某个文件中包含多个文本片段，用空格分隔，现在需要查找包含用户输入内容的片段：用户输入 *cat*，查找包含 *cat* 的行，直接的思路是 *^.\*cat.\*\$*（使用多行模式，同时不能指定单行模式），假设用户输入的内容保存在变量 *userInput* 中，就应当用 *"^.\*" + userInput + ".\*"* 得到查找用的正则表达式。

如果用户输入 *cat*，这么做当然没有问题。如果用户输入的是 *ca\*t*，得到的正则表达式就是 *^.\*ca\*t.\*\$*：本意是查找包含字符串 *ca\*t* 的行，但 *\** 是正则表达式中的元字符，正则表达式 *ca\*t* 能匹配字符串 *cat*、*caat*、*caaat*，却不能匹配 *ca\*t*，就会产生错误。更麻烦的是，恶意用户可能会输入 *a(b\*b\*)b\** 之类的字符串，这样的正则表达式匹配起来会消耗非常多的资源，这就是 **正则表达式拒绝服务攻击**（☞ 135），严重时可能导致服务瘫痪。

为解决这类问题，一些语言中的正则表达式提供了特殊结构，彻底消除元字符的特殊含义，提供真正“安全”的表达式（也就是普通字符串）。最常见的做法是在需要消除元字符特殊含义的正则表达式两端添加 *\Q* 和 *\E*，也就是 *\Q... \E*（其中 *Q* 表示“引用文本 Quoted”，而 *E* 表示“引用文本结束 End”），像例 6-7 那样。此时正则表达式 *\Qca\*t\E* 就不再能匹配 *cat*、*caat*、*caaat*，只能匹配 *ca\*t*。*\Q... \E* 只对在它内部的子表达式生效，所以在正则表达式 *^.\*\Qca\*t\E.\*\$* 中，*^.\*\$* 和 *.\*\$* 会按照正则表达式的默认规则匹配，只有中间的 *ca\*t* 作为普通字符串匹配，整个表达式能够匹配的就是确定包含 *ca\*t* 的行。

#### 例 6-7 用 *\Q... \E* 消除元字符的特殊含义

```
//因为 Python 不支持 \Q... \E 的表示法，现在以 Java 语言举例
//没有使用 \Q... \E
"cat".matches("ca*t");      // => True
"caat".matches("ca*t");     // => True
```



```
"caaat".matches("ca*t");    // => True
"ca*t".matches("ca*t");    // => False
//使用\Q...\E，注意在正则文字中必须双写反斜线为\\Q...\E
"cat".matches("\\Qca*t\\E");    // => False
"caat".matches("\\Qca*t\\E ");    // => False
"caaat".matches("\\Qca*t\\E ");    // => False
"ca*t".matches("\\Qca*t\\E ");    // => True
//对正则表达式中的某个子表达式使用\Q...\E
"regex like ca*t and so on".matches("^.*\\Qca*t\\E.*$");    // => True
"has one cat".matches("^.*\\Qca*t\\E.*$");    // => False
```

`\Q...\E`并不是所有语言都支持的，在本书介绍的语言中，明确支持它的有 Java 和 PHP。其他语言虽然不支持`\Q...\E`的记法，但其中一些也提供了专门的处理函数消除元字符的特殊含义；比如.NET 的 `Regex.Escape(text)`，它接收一个字符串，返回的字符串中的所有元字符都做了转义处理，消除了特殊含义，从其结果字符串生成的正则表达式，就只能匹配与 `text` 完全一样的字符串，代码见例 6-8。

例 6-8 用专门的函数消除元字符的特殊含义

```
Regex.IsMatch("ca*t", Regex.Escape("ca*t"));    // => True
Regex.IsMatch("(ab)", Regex.Escape("(ab)"));    // => True
```

表 6-8 列出了各种语言中消除元字符特殊含义的函数。

表 6-8 常用语言中消除元字符特殊含义的函数

语言	函数
.NET	<code>Regex.Escape(text)</code>
Java	<code>Pattern.quote(text)</code>
PHP	<code>preg_quote(text)</code>
Python	<code>re.escape(text)</code>
Ruby	<code>Regexp.quote(text)</code> <code>Regexp.escape(text)</code>

注：PHP 中的正则表达式两端必须出现分隔符，所以 `preg_quote()`可以设定第二个参数明确指定分隔符；否则，假如分隔符是`/`，而 `text` 又包含`/`，就可能发生冲突。

请记住，使用正则表达式时，仔细分辨并消除用户输入字符串中元字符的特殊含义，是不可忽略的步骤。

### 6.1.4 字符组中的转义

在正则表达式中，如果需要表示作为元字符的普通字符（比如`*`、`?`、`(`等），就需要使用转义，这一点不存在疑义。特殊的是，常见的元字符出现在字符组内部基本都不算元字符，也就是说，它们在字符组内部出现时，不需要转义。代码如例 6-9 所示。

例 6-9 字符组内部几乎没有元字符

```
re.search(r"[*]", "*") != None    # => True
re.search(r"[?]", "?") != None    # => True
re.search(r"[()]", "(") != None    # => True
```

之前讲过，字符组有自己的元字符规定，也有相应的转义规定：在字符组内部，只有三个字符需要转义。一个是闭方括号`]`，如果不是作为字符组结束标志的闭方括号，则必须写成`\]`，比如`[0\]9]`，它可以匹配的字符是`0`、`]`、`9`；一个是横线`-`，如果不是用于范围表示法（比如`[0-9]`），必须写成`\-`，比如`[0\ -9]`，它可以匹配的字符是`0`、`-`、`9`，当然如果它紧跟在开方括号之后，也可以不用转义，`[0\ -9]`和`[-09]`是等价的（我更推荐后一种写法，因为更清晰简洁）；还有一个需要转义的字符是`^`，如果它不是用于排除型字符组（`[^ab]`），则应当写成`\^`，比如`[\^ab]`，它可以匹配除`^`、`a`、`b`之外的任何字符，如果它不是紧跟在开方括号之后，也不用转义，`[\^ab]`、`[a^b]`、`[ab^]`是完全等价的（我更推荐后两种写法，因为更清晰简洁）。从例 6-10 可以看到，这两种写法的结果是相同的。

例 6-10 字符组中三个需要转义的元字符

```
#转义的]
re.search(r"[0\]9]", "]") != None    # => True
#未转义的-
re.search(r"[0-9]", "-") != None      # => False
#转义的-
re.search(r"[0\ -9]", "-") != None    # => True
re.search(r"[-09]", "-") != None      # => True
#未转义的^
re.search(r"[^ab]", "^") != None      # => True
re.search(r"[^ab]", "a") != None      # => False
#转义的^
re.search(r"[\^ab]", "^") != None      # => True
re.search(r"[\^ab]", "a") != None      # => True
re.search(r"[a^b]", "^") != None      # => True
re.search(r"[a^b]", "a") != None      # => True
```

## 6.2 正则表达式的处理形式

之前讲解正则表达式时主要使用的 Python 语言中，其具体办法是调用 `re` 这个 package 中的函

数（方法），比如 `re.search()`、`re.findall()`、`re.sub()`。选择合适的函数，将正则表达式和字符串传入，这是一种常见的方式，比如 PHP 也是如此。但是，正则表达式并不只有这一种处理形式，比如 Java 语言的处理形式就与 Python 和 PHP 大不一样，所以下面简单介绍常见的处理形式。

### 6.2.1 函数式处理

在函数式处理中，正则表达式的常见操作（查找、替换等）都有对应的函数，执行这些操作时，调用对应的函数，将正则表达式和字符串作为参数传入即可，Python、PHP 是函数式处理的典型代表。

#### Python

```
#查找
re.findall(r"\d+", "123 45 6")
#替换
re.sub(r"\d+", r"[\g<0>]", " 123 45 6")
```

#### PHP

```
//查找
preg_match("/\d+/", "123 45 6");
//替换
preg_replace("\d+", "[$0]", " 123 45 6");
```

### 6.2.2 面向对象式处理

Java 和 .NET 之类的语言中的正则表达式采取了不同的处理形式：进行正则表达式处理之前，必须生成专门的正则表达式对象（在不同的语言里，对象所属的类名不同），再调用此对象的成员函数。

#### Java

```
//必须先生成 Pattern 对象，再生成 Matcher 对象
Pattern pattern = Pattern.compile("\\d+");
Matcher matcher = pattern.matcher("123 45 6");
while(matcher.find()) {
    //other code here
}
```

#### .NET

```
//必须先生成 Regex 对象，再生成 Match 对象
Regex regex = new Regex(@"\d+");
Match match = regex.Match("123 45 6");
while (match.Success) {
    //other code here
    match = match.NextMatch();
}
```

## 6.2.3 比较

同样是进行查找操作，如果使用函数式处理，只需要调用对应的函数；但使用面向对象式处理，需要逐步生成各种对象，再调用对象自身的方法。看起来，后者要麻烦很多，这种处理方式有什么好处呢？

之前讲过，正则表达式并不等于字符串，即便正则表达式是以字符串形式给出的，进行正则表达式操作之前，必须首先生成专用的“正则表达式对象”。函数式处理隐去了生成的过程，感觉更加直接；面向对象式处理则暴露了生成的过程，感觉更加细致。

只要执行正则表达式操作，就会产生正则表达式对象；所以，面向对象式处理的代码虽然更烦琐，但是如果需要重复用到某些表达式，它的效率往往远高于函数式处理——在面向对象式处理中，可以将已经生成的对象作为变量保存起来，就不必重复生成了。

面向对象式处理的步骤分明、效率更高，而函数式处理的好处是方便顺手、代码简洁；最好的办法是，根据应用场合的不同，采取不同的处理方式：如果正则表达式只是单次使用，则选择函数式处理；如果正则表达式需要重复使用，则选择面向对象式处理。实际上，许多编程语言也提供了对应的设计，比如 Python 提供了一些面向对象式处理的方法，Java 和 .NET 也提供了一些函数式处理的方法。

在 Python 中进行面向对象式处理，可以先调用 `re.compile()` 生成专门的 `RegexObject` 对象（也就是正则表达式对象），在进行正则表达式操作时，可以把这个对象作为参数，传给对应的函数；另一方面，`re` 中各个函数也可作为对象自身的成员方法，所以调用生成好的 `RegexObject` 对象的成员方法，也是可以的，如例 6-11 所示。

例 6-11 Python 中可以使用两种处理方式

```
#函数式
re.findall(r"\d+","123 45 6")
#面向对象式
#生成 re 对象
digitRegex = re.compile(r"\d+")
#调用 re.findall 方法
re.findall(digitRegex, "123 45 6")
#或直接调用 re 对象的方法
digitRegex.findall("123 45 6")
```

Java 语言虽然采用面向对象式处理，但也提供了一些便捷函数，比如 `Pattern.matches(regex, input)`，它用来验证 `input` 能否由 `regex` 匹配，如例 6-12 所示。

例 6-12 Java 中也可以使用函数式处理

```
Pattern.matches("\\d+", "123 45 6");
```

```
//它等价于
Pattern.compile("\\d+").matcher("123 45 6").matches();
```

.NET 中也提供了类似的方法，比如 `Regex.IsMatch(String regex, String input)`，它也用来验证 `input` 能否由 `regex` 匹配，代码如例 6-13 所示。

例 6-13 .NET(C#)中也可以使用函数式处理

```
Regex.IsMatch(@"\d+", "123 45 6");
//它等价于
(new Regex(@"\d+").Match("123 45 6")).Success;
```

## 6.2.4 线程安全性

如果采用面向对象式处理，可以将生成的正则表达式对象存储在变量中反复使用，以提高效率。在单线程编程的环境下，“反复使用”当然没有问题，如果多个线程共用同一个正则表达式对象，结果会如何呢？

仔细观察面向对象式处理可以发现，在面向对象式处理中，一般会出现两个对象：一个是单独对应正则表达式的，可以叫做“正则表达式对象”，比如 Java 中的 `Pattern`、.NET 中的 `Regex`、Python 中的 `RegexObject`；另一个是在正则表达式与字符串联系起来时生成的对象，可以叫做“匹配结果对象”，比如 Java 中的 `Matcher`、.NET 中的 `Match`、Python 中的 `MatchObject`<sup>1</sup>。这几个对象之间的关系如图 6-2 所示。

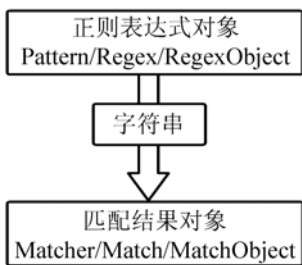


图 6-2 正则表达式操作中对象间的关系

正则表达式对象本身基本没有什么状态可言，所以这个对象总是线程安全的，可以由多个线程共享。匹配结果对象则需要维护自身状态，比如当前匹配是否成功，当前匹配的开始或结束偏移值等。如果多个线程共享同一个匹配结果对象，有可能遇到这样的情况：某个线程正在查询当前匹配的开始位置，另一个线程调用了 `matcher.find()`（Java 中）或者 `Match.nextMatch()`（.NET

<sup>1</sup> 也有例外，比如 JavaScript 中的 `RegExp`，它既是正则表达式对象，其中又保存了与具体匹配相关的信息。不过 JavaScript 一般不会使用多线程编程，所以这个例外不做详细介绍。

中)之类的方法,尝试进行下一次匹配,就可能产生混乱。所以匹配结果对象一般不是线程安全的,也不应由多个线程共享。

所以,最理想的办法是:多个线程可以共享同一个正则表达式对象,节省时间;但操作不同文本时,应当针对各个线程生成专属的匹配结果对象。

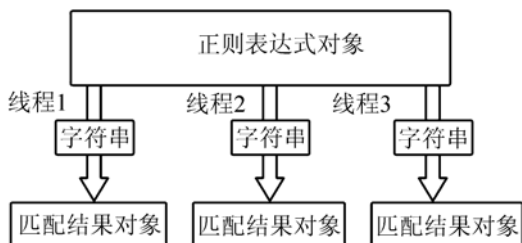


图 6-3 多线程环境下正确的处理模式

例 6-14 多线程环境下错误使用正则表达式 (以 Java 为例)

```
//多个线程共享同一个 Matcher
Pattern pattern = Pattern.compile("\\d+");
Matcher matcher = pattern.matches("123 45 6");
...
//thread 1
while(matcher.find()) {
    //other code here
}
...
//thread 2
matcher.reset("78 9 0");
while(matcher.find()) {
    //other code here
}
```

例 6-15 多线程环境下正确使用正则表达式 (以 Java 为例)

```
//各个线程共享同一个 Pattern
Pattern pattern = Pattern.compile("\\d+");
...
//thread 1
Matcher matcher = pattern.matches("123 45 6");
while(matcher.find()) {
    //other code here
}
```

```
...
//thread 2
Matcher matcher = pattern.matches("78 9 0");
while(matcher.find()) {
    //other code here
}
```

6.3 表达式中的优先级

正则表达式千变万化，都是由之前介绍的字符组、括号、量词等等基本结构组合而成的，只要掌握了组合的规则，面对再复杂的表达式，都能把结构梳理清楚。

仔细观察会发现，正则表达式的元素之间的组合关系只有 4 种。

普通拼接	abc
括号	(abc)
量词	a*b
多选结构	ab cd

注：“普通拼接”可能是最常见的组合关系，只是平时不一定明确意识到，正则表达式 abc 就是 a 和 b 和 c 的普通拼接。a(bc)则可以看做 a 和(bc)的拼接。

除列出的 4 种组合关系外，正则表达式中的其他结构，比如环视( ?=... )，都可以视为单独的元素，其中的正则表达式最终可以归类到上面的 4 种组合关系中。

所以，真正需要关心的就是这 4 种组合的优先级，详见表 6-9。

表 6-9 正则表达式中的优先级

优先级	组合	说明
1	(regex)	整个括号内的子表达式成为单个元素
2	* ? +	限定之前紧邻的元素
3	abc	普通拼接，元素相继出现
4	a bc	多选结构

注：数字越小，优先级越高。

有了表 6-9，就可以拆解各种表达式了，表 6-10 给出了几个简单的例子。

表 6-10 正则表达式中的优先级举例

正则表达式	能匹配的字符串	说明
ab	ab	普通拼接
ab+	ab abb abbb	量词的优先级高于普通拼接
(ab)+	abc aabc aabc	括号的优先级高于量词
ab cd	ab cd	多选结构的优先级低于量词
(ab c)d	abd cd	括号的优先级高于多选结构和量词

如果你觉得上面这几个例子都不难理解，来看个更复杂的表达式 `ab*(cd|e+)?|fg`，它的优先级划分关系如图 6-4 所示。

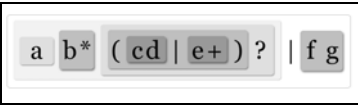


图 6-4 表达式的优先级

注：颜色越深，优先级越高。

判断优先级并不困难，一般来说，根据上面介绍的知识逐步拆解即可，但是多选结构值得单独强调。

常见的多选结构是以竖线配合括号 `(option1|option2)` 的形式给出来的；但是，多选结构并非必须要与括号连用，如果没有括号，单独出现竖线 `|`，也可以实现多选结构的功能，比如正则表达式 `ab|cd`，它能匹配的字符串和 `(ab|cd)` 是完全一样的。

比较起来，后一种写法更加清楚，因为两端的括号明确限定了竖线所标识的多选结构的范围，所以很难混淆 `^(ab|cd)$` 和 `^ab|cd$`，后者等于 `(^ab|cd$)`；而且在匹配之后，还可以使用分组及反向引用等功能（`\1`、`\2` 之类），明确得到匹配的文本（不用括号则做不到这点），进行下一步操作。但是，分组又会降低效率，尤其是在正则表达式非常复杂，或者要处理的文本非常多时，使用分组可能会严重影响效率。

考虑到这些因素，我推荐的做法是：应用多选结构一定不能省略括号，但同时，为了避免分组影响效率，如果不需要提取捕获文本，应当把普通的括号 `(...)` 改为非捕获型括号 `(?:...)`，这样既明确标识了多选结构，又免去了分组捕获的成本。不过在本书中，为了格式简单清晰考虑，除非特别说明，一般不使用非捕获型括号。





## 第二部分

## 第7章 Unicode

### 7.1 关于编码

通常，英文编码较为统一，都采用 ASCII 编码<sup>1</sup>或可以兼容 ASCII 编码（即编码表的前 127 位与 ASCII 编码一致，常见的各种编码，包括 Unicode 编码都是如此）。也就是说，英文字母、阿拉伯数字、英文的各种符号，在不同编码下的码值（Code Point）基本是一样的<sup>2</sup>，比如字母 A，其码值总是 41；中文的情况则不同，常见的中文编码有 GB18030（也就是 CP54936，主要是在 Windows 平台下使用。早期是 GBK<sup>3</sup>，也就是 CP936，如今采用的 GB18030 与 GBK 是兼容的，考虑到大家习惯说“GBK 编码”，下文也采用“GBK 编码”的说法）和 Unicode（主要用于 Linux/UNIX、Mac OS）两种，同一个中文字符在两种编码下的码值并不相同。比如“发”，在 GBK 编码下的码值为 b7 a2，而在 Unicode 下其码值为 53 d1（根据惯例，Unicode 字符一般记为 \uhex 的形式，其中 hex 是码值的十六进制表示，所以“发”的 Unicode 记法就是 \u53d1）。

为方便下面的讲解，这里先约定两种提法：

**ASCII 字符**，即 ASCII 编码表中的字符（也就是码值在 0~127 之间的字符，不包括扩展 ASCII 字符），每个字符用一个字节表示。常见的英文字符和半角标点符号，都属于 ASCII 字符。

**非 ASCII 字符**，即 ASCII 编码表之外的字符，在本书中指多字节字符。中文字符属于“非 ASCII 字符”，它们在 GBK 编码中一般占用 2 个字节<sup>4</sup>，在 UTF-8 编码中占用 3 个字节。常见的 Unicode 编码是 UTF-8，它是 Unicode 的一种存储格式，UTF-16 也是一种存储格式。无论是 UTF-8 还是 UTF-16，同一字符的 Unicode 编码是不变的，只是存储时的表现形式不同。比如“发”的码值是 53 d1，但以 UTF-8 编码保存时占用 3 个字节，分别是 e5 8f 91，保存为 UTF-16 编码时占用 2 个字节，分别是 53 d1。

---

<sup>1</sup> 严格说起来应该是“Unicode 字符集”而不是“Unicode 编码”，因为同一种 Unicode 字符在不同的编码（或者叫转换格式，比如 UTF-8、UTF-16）中有不同的码值。

<sup>2</sup> 也有例外，比如 EBCDIC 编码就是与 ASCII 不兼容的单字节编码，但我们遇到的机会很少。

<sup>3</sup> “真正的”国标是 GB2312 编码，GBK 编码是对它的扩充，Windows 系列操作系统采用 GBK 编码。

<sup>4</sup> 严格来说，按照 GB18030 编码，可能出现 4 个字节的字符，但绝大多数中文字符都占用 2 个字节。

## 7.2 推荐使用 Unicode 编码

常见的正则表达式的文档都是关于英文（ASCII 字符）的，英文开发者通常也只需要处理 ASCII 字符，不需要处理中文这类多字节的字符。不过，依照处理 ASCII 字符的方式处理中文字符，就有可能出错。

举个例子，许多文档说点号 `.` 可以匹配“除换行符 `\n` 之外的任意字符”，但这可能只适用于单字节字符，因为点号匹配的其实只是“除换行符 `\n` 之外的任意字节”而已。

```
Python 2
re.search(r"^.$","发") != None # =>False
PHP 4.x/5.x
preg_match("/^.$/", "发"); // =>0, 表示成功匹配的次数为 0
Ruby 1.8
"发" =~ /^.$/ != nil # => False
```

之所以会出现这种情况，是因为正则表达式在匹配时“看到”的字符串是多个字节（在 UTF-8 编码环境<sup>1</sup>下占用 3 个字节；在 GBK 编码环境下占用 2 个字节），于是将它们全部视为单字节字符，而没有理解这几个字节表示的是“发”；因为点号 `.` 只能匹配一个字符，结果当然不能正确匹配。要解决这个问题，必须显式指定编码，让正则表达式处理程序正确识别多字节字符，Python 中的做法是在字符串开头的引号之前写上 `u` 字符，如例 7-1 中的 Python 代码所示。

### 例 7-1 显式指定 Unicode 的字符串（Python 2）

```
#单个“发”字符被拆散为 3 个单字节字符，码值分别为 e5 8f 91
print repr(r'发')
'\xe5\x8f\x91'

print repr(u'发')
u'\u53d1'
```

如果使用的是 Python 2、Ruby 1.8 或 PHP，都需要显式指定 Unicode<sup>2</sup>。而在 Python 3、Java、.NET 和 Ruby 1.9 中，字符串默认采用 Unicode 编码，所以不存在上面的问题。

```
Python 2
re.search(r"^.$", u"发") != None # => True
```

---

<sup>1</sup> “编码环境”指的是源代码和它所处理的数据文本采用的编码，“UTF-8 编码环境”的意思是，源代码和处理的数据都采用 UTF-8 编码，而“GBK 编码环境”的意思是，源代码和处理的数据都采用 GBK 编码。

<sup>2</sup> 如果在 Python 2 中使用 Unicode 字符串，必须在源代码中显式指定所使用的编码，具体做法是，在源代码的顶端写上这样一行：`# -*- coding: UTF-8 -*-`。如果在 Ruby 1.8 中，则必须在源代码的顶端写上：`# encoding:UTF-8`。

```

PHP 4.x/5.x
preg_match("/^.$$/u", "发") != 0;           // => true
Ruby 1.8
"发" =~ /^.$$/u                             # => 0

```

细心的读者会发现，PHP 和 Ruby 是对正则表达式指定 Unicode，Python 2 却是对字符串指定 Unicode。如果熟悉 Python 文档，就会知道那里面同时说明了，正则表达式也可以指定 Unicode 模式，这是怎么回事呢？

一般来说，正则表达式的操作可以简要概括为“用正则表达式匹配字符串”。在 Python 2 中，字符串默认采用 ASCII 编码，多字节字符就会拆分为多个单字节字符，比如“发”字被视为由多个单字节字符构成的字符串，点号 `.` 只能匹配其中的一个单字节字符。如果将整个字符串显式指定为 Unicode 字符串（也就是 `u"发"`），则“发”字会被视为单个字符，点号 `.` 可以匹配。

同时，因为正则表达式也是以字符串的形式提供的，所以如果在正则表达式的字符串中使用中文字符，这个字符串也应该设定为 Unicode 字符串；否则就可能出错，尤其在字符组中更是如此。之前已经讲过，`r"^[发]$"` 之类的原生字符串更适合作正则表达式，如果要同时指定这个字符串为 Unicode 字符串，则应当写作 `ur"^[发]$"`，而不是 `ru"^[发]$"`，如例 7-2 所示，请注意 `u` 和 `r` 的顺序。

#### 例 7-2 作为正则表达式的字符串也应当指定 Unicode

```

re.search(r"^[发]$", u"发") != None        # => False
re.search(ur"^[发]$", u"发") != None       # => True

```

如果没有指定 `^[发]$"` 为 Unicode 编码，其中的“发”字同样会被拆散为多个字节，所以这个字符组的内容并不是字符“发”，而是“发”所对应的多个单字节字符。为说明这一点，下面使用 Python 的 `re` 中定义的 `re.DEBUG` 这个常量，通过它可以看到正则表达式的内部结构。如例 7-3 所示（程序在 UTF-8 编码下运行），只有在明确指定正则表达式的字符串也使用 Unicode 编码的情况下，才能正确识别字符组内部的“发”。

#### 例 7-3 通过 `re.DEBUG` 观察正则表达式内部的结构

```

re.compile(r"[发]", re.DEBUG)
in
    literal 229
    literal 143
    literal 145

re.compile(ur"[发]", re.DEBUG)
literal 21457

```

在 PHP 和 Ruby 1.8 中并不存在“Unicode 字符串”，所以无法修改字符串的属性来指定 Unicode 编码。不过如果将正则表达式指定为 Unicode 模式，正则表达式处理程序就可以正确识别字符串中的 Unicode 字符。所以如果用 PHP 或 Ruby 的正则表达式处理 Unicode 字符串，一定不要忘记指定 Unicode 模式。表 7-1 总结了常用语言中点号对多字节字符的匹配情况。

表 7-1 点号对 Unicode 字符的匹配（采用 UTF-8 编码）

字符串	正则表达式	语言	是否显式指定 Unicode 模式	可否匹配
发	^.\$	.NET	否	可
	^.\$	Java	否（无须指定）	可
	^.\$	JavaScript	否（无法指定）	结果取决于浏览器
	/^.\$/	PHP	否	否
	/^.\$/u	PHP	是	可
	^.\$	Python 2	否	否
	^.\$	Python 3	否	可
	/^.\$/	Ruby 1.8	否	否
	/^.\$/u	Ruby 1.8	是	可
	/^.\$/	Ruby 1.9	否	可

注 1：PHP 和 Ruby 的正则表达式本身不包含分隔符（分隔符可以有很多种，常见的是反斜线/），但 PHP 指定 Unicode 模式必须在后一个分隔符之后写 u，所以在这里将分隔符也写出来。

注 2：Python 2 中虽然提供了 Unicode 模式，但即便采用这种模式，点号仍然不能匹配多字节字符，除非显式指定字符串为 Unicode 字符串（在 Python 2 中，默认是 ASCII 字符串）。

注 3：Ruby 1.9 在不同平台下的表现不一样，经过作者测试，在 Linux 下采用 Unicode 匹配规则，而 Mac 和 Windows 的情况通常是采用 ASCII 匹配规则。

现在再来解释，虽然 GBK 编码也是多字节编码，为什么不推荐使用：常见的中文字符编码有 GBK 和 Unicode 两种，GBK 编码使用确实也很多（尤其是在 Windows 平台上），但如果需要使用正则表达式处理中文，我强烈推荐使用 Unicode 字符。这样选择的最重要原因是，正则表达式处理程序一般只能准确识别 Unicode 字符的边界。

举例来说：要求匹配“收”（GBK 编码是 ca d5）或者“发”（GBK 编码是 b7 a2），很自然会想到使用字符组[收发]。这个思路是对的，但从例 7-4 可以看到，即便采用 GBK 编码，正则表达式处理程序仍然不能识别出“2 个多字节字符构成的字符组”，而将其视为“4 个单字节字符构成的字符组”，其值分别为 202、213、183、162，换算成十六进制，正好是“收”和“发”对应的 ca、d5 和 b7、a2。用它来匹配没有指定 Unicode 编码的字符“罚”（它的 GBK 编码是 b7 a3，也就是十进制的 183 163），就会产生错误——虽然“罚”字既不等于“收”也不等于“发”，但码

值为 b7 的这个单字节字符可以匹配！

#### 例 7-4 GBK 编码环境下的字符组会出现错误匹配

```
re.compile(r"[收发]", re.DEBUG)
in
    literal 202
    literal 213
    literal 183
    literal 162
print repr("罚")
'\xb7\xa3'
re.search(r"[收发]", "罚") != None      # =>True
```

如果使用 Unicode 编码环境，正则表达式处理程序可以准确识别字符的边界，所以知道[收发]包含码值为 65 36 和 53 d1 的两个多字节字符，而“罚”的码值是 7f 5a，因此必然不会发生错误匹配。代码如例 7-5 所示。

#### 例 7-5 Unicode 编码环境下的字符组可以避免错误匹配

```
#十六进制的 6536 和 53d1 分别等于十进制的 25910 和 21457
re.compile(ur'[收发]', re.DEBUG)
in
    literal 25910
    literal 21457
re.search(ur"[收发]", u"罚") != None      # =>False
```

你或许会想，如果使用 GBK 编码环境，但仍然指定字符串为 Unicode 模式，是否可以解决问题？例 7-6 说明，在 GBK 编码环境下指定字符串为 Unicode 模式，码值会发生错乱，“发”字的 Unicode 编码分明是 53 d1，GBK 编码环境下虽然可以指定 Unicode 模式，码值却仍然是 GBK 的 b7 a2。

#### 例 7-6 GBK 编码环境下使用 Unicode 字符串会出错

```
#Unicode 编码环境
print repr(u'发')
u'\u53d1'

#GBK 编码环境
print repr(u'发')
u'\xb7\xa2'
```

如果因为某些限制只能使用 GBK 编码，有一个窍门准确保证 `[收发]` 的匹配，就是把字符组 `[收发]` 改成多选结构 `(收|发)`。从例 7-7 可以看到，此时如果要匹配成功，只能是两个连续的字节 `ca d5` 或者 `b7 a2`，而“罚”字的两个字节为 `b7 a3`，无法匹配。

#### 例 7-7 GBK 编码的字符组应改为多选结构避免错误匹配

```
re.search(r"(收|发)", "罚") != None      # => False
re.search(r"(收|发)", "收") != None      # => True
re.search(r"(收|发)", "发") != None      # => True
```

即使这样，排除型字符组还是可能匹配出错。如果希望匹配一个“收”和“罚”之外的字符，自然的思路就是使用排除型字符组 `[^收发]`。但是根据上面的讲解，这样“排除”的并不是 2 个多字节字符，而是 4 个单字节字符，其码值分别为 `ca d5 b7 a2`。另一方面，“罚”字的 GBK 编码为 `b7 a3`，`b7` 这个字节被“排除”了，所以结果如例 7-8 所示，“罚”字不能由 `[^收发]` 匹配，Unicode 编码环境则没有这样的问题。

#### 例 7-8 排除型字符组的错误匹配

```
#GBK 编码环境
re.search(r"[^收发]", "罚") != None      # => True
#Unicode 编码环境
re.search(ur"[^收发]", u"罚") != None    # => False
```

所以，为避免出现错误匹配，总的原则就是，如果使用正则表达式处理含有多字节字符的文本，则尽量使用 Unicode 编码环境。

## 7.3 Unicode 匹配规则

上一节提到，Python 文档说明，除了指定字符串为 Unicode 字符串，还可以指定正则表达式使用 Unicode 模式，这又是怎么回事呢？

不妨回头仔细想想读过的文档，正则表达式中的 `\d` 和 `\w` 都是如何解释的？许多人的第一反应是：`\d` 等价于 `[0-9]`，`\w` 等价于 `[0-9a-zA-Z_]`，因为有些文档说明了这种等价关系。不过也有些文档说：`\d` 匹配数字字符，`\w` 匹配单词字符，`\s` 匹配空白字符；但是没有说明数字字符、单词字符、空白字符到底是哪些字符。

一般来说，数字字符就是 `[0-9]`，单词字符就是 `[0-9a-zA-Z_]`，空白字符则包括空格、回车等字符，但这只是 ASCII 编码中的情况，在 Unicode 编码中并非如此。

因为包括了多种语言和字符，所以在 Unicode 编码中，全角数字 0、1、2 之类也算作“数字字符”，可以由 `\d` 匹配；中文字符，也可以算作“单词字符”，由 `\w` 匹配；同样的道理，中文



的全角空格（码值为 30 ff），也可以算作“空白字符”，由 `\s` 匹配。所以，如果在 Python 2 中指定了正则表达式使用 Unicode 模式（最简单的方式就是在正则表达式的开头指定模式修饰符 `(?u)`），`\d`、`\w`、`\s` 就能匹配全角数字、中文字符、全角空格。对于这种情况，本书中称为 **Unicode 匹配规则**；相应地，之前 ASCII 编码中的匹配，称为 **ASCII 匹配规则**。

例 7-9 排除型字符组的错误匹配

```
#字符均为全角
re.search(r"^\d$", u" 1 ") != None      # => False
re.search(r"^\w$", u"发") != None       # => False
re.search(r"^\s$", u" ") != None        # => False

re.search(r"(?u)^\d$", u" 1 ") != None   # => True
re.search(r"(?u)^\w$", u"发") != None   # => True
re.search(r"(?u)^\s$", u" ") != None    # => True
```

表 7-2 详细介绍了两种匹配规则下 `\d`、`\w`、`\s` 的匹配。相应地，此时 `\D` 匹配 `\d` 不能匹配的字符，`\W` 匹配 `\w` 不能匹配的字符，`\S` 匹配 `\s` 不能匹配的字符。

表 7-2 Unicode 模式下 `\d`、`\w`、`\s` 的匹配

简记法	ASCII 匹配规则	Unicode 匹配规则
<code>\w</code>	<code>[0-9a-zA-Z_]</code>	<code>[\p{Ll}\p{Lu}\p{Lt}\p{Lo}\p{Lm}\p{Nd}\p{Pc}]</code>
<code>\d</code>	<code>[0-9]</code>	<code>[\p{Nd}]</code>
<code>\s</code>	<code>[\t, \n, \v, \f, \r, \x20]</code>	<code>[\f\n\r\t\v\x85\p{Z}]</code>

注：`\p{L}`表示任意语言中的字母字符（包括英文字母和汉字）；  
`\p{M}`表示用来与其他字符结合的字符（声调、元音变化符等）；  
`\p{Nd}`表示任何书写系统中的 0~9 的字符，汉字全角字符 1、2 等也算；  
`\p{NI}`表示形如字符的数字，比如罗马数字；  
`\p{Pc}`表示类似下画线之类的标点字符；  
`\p{InEnclosedAlphanumerics}`表示被包围的数字或字符，比如①；  
`\p{So}`表示数字符号、货币符号、组合字符之外的符号字符。  
 具体细节请参考本章末 Unicode Property 和 Unicode Block 的详细列表。

有时候，这样的规定确实让人抓狂：假设你希望用正则表达式 `\d{6,12}` 来验证一个长度在 6 到 12 之间的数字字符串，却没留意 `\d` 能匹配全角数字，验证就可能出错；所以，一定要注意此类问题，表 7-3 总结了常用语言中的匹配规则。

表 7-3 常用语言中的匹配规则

语言	字符组简记法的匹配规则
.NET	默认采用 Unicode 匹配规则，但可以显式指定采用 ASCII 匹配规则
Java	ASCII 匹配规则
JavaScript	ASCII 匹配规则
PHP	ASCII 匹配规则
Python 2	默认采用 ASCII 匹配规则，但可以显式指定采用 Unicode 匹配规则
Python 3	默认采用 Unicode 匹配规则，但可以显式指定采用 ASCII 匹配规则
Ruby 1.8	默认采用 ASCII 匹配规则，显式指定 Unicode 模式之后，只有\w 使用 Unicode 匹配规则
Ruby 1.9	Unicode 匹配规则

注 1：一般来说，单词边界\b 能匹配的位置是：一端是单词字符，一端不是单词字符（也可以什么都没有），其中单词字符的规定与\w 一样。Java 中则不是这样，它的细节比较复杂，具体请参考 Java 的章节。

注 2：在 Python 3 中可以在表达式最开始用(?a)指定 ASCII 模式。

注 3：本书中说的 Ruby 1.9 指的是 Ruby 1.9.1 及以上版本，并不包含 Ruby 1.9.0，因为这个版本使用非常少，而且有一些非常奇怪的表现。

## 7.4 单词边界

单词边界之前已经介绍过（☞ 59），它的准确解释是：一端必须出现\uw 能匹配的字符，另一端不出现\uw 能匹配的字符。在 JavaScript、PHP、Python 2、Ruby 中，\w 只能匹配[0-9a-zA-Z\_]。所以在这些语言中，\b\w+\b 能用来匹配几乎所有的英文单词，在例 7-10 中用它来提取一段文本中的所有单词。

例 7-10 用单词边界提取单词

```
re.findall(r"\b\w+\b", "a sentence contains a lot of words")
["a", "sentence", "contains", "a", "lot", "of", "words"]
```

之所以说“几乎所有”英文单词，是因为有些情况更加复杂：假设字符串中包含中文，之间用空格字符分隔，比如“他们 我们去见他们的时候……”。如果希望匹配单独出现的“他们”，最自然的想法似乎是使用正则表达式\b 他们\b，结果却是不对的。

例 7-11 用单词边界提取中文单词会出错

```
re.findall(r"\b 他们\b", "他们 我们去见他们的时候……")
[]
```

原因在于，默认情况下 `\w` 是按照 ASCII 匹配规则的。“他”字之前的位置，左侧没有字符，不能由 `\w` 匹配，右侧的“他”字也不能由 `\w` 匹配；“们”字右侧的情况也是如此。所以，`\b 他们\b` 无法匹配“他们”。要解决这个问题可以采用 Unicode 匹配规则，此时 `\w` 就不只能匹配 `[0-9a-zA-Z_]`，还能匹配其他语言中的单词字符（包括中文字符），程序代码如例 7-12 所示。

例 7-12 在 Unicode 模式下提取中文单词

```
re.findall(r"(\u)\b 他们\b", "他们 我们去见他们的时候……")
['他们']
```

但是，这样又出现了新的问题。比如字符串“学习 regex”，其中的 `regex` 在非 Unicode 匹配规则下是可以由 `\b regex\b` 匹配的；在 Unicode 匹配规则下却行不通，因为此时 `\w` 既能匹配中文字符又能匹配英文字符，`\b regex\b` 反而无法匹配 `regex` 了，代码见例 7-13。

例 7-13 在 Unicode 模式下提取单词仍然有问题

```
re.findall(r"\b regex\b", "学习 regex")
["regex"]

re.findall(r"(\u)\b regex\b", "学习 regex")
[]
```

为了形象说明在 ASCII 匹配规则和 Unicode 匹配规则下 `\b` 能够匹配的“单词边界”，（因为通常处理的文本包含中文和英文，所以只是简要列出了各分类的常用字符），下面用图 7-1 和图 7-2 分别说明。其中深色区域是必须出现的，而浅色的区域是可以出现、也可以不出现的（比如文本首尾位置的空字符串）。

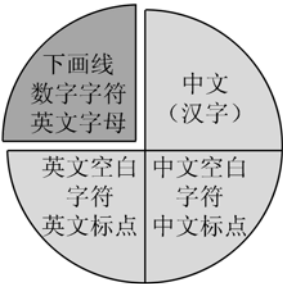


图 7-1 ASCII 匹配规则下的单词边界



图 7-2 Unicode 匹配规则下的单词边界

表 7-4 举列列出了在 Unicode 匹配模式下，`\b` 匹配的几种典型的场景。

表 7-4 在 Unicode 匹配模式下，单词边界的几种典型场景

字符串	正则表达式	能否匹配	说明
tom, jerry（半角标点）	\btom\b	能	\w 不能匹配半角标点
tom, jerry（全角标点）	\btom\b	否	\w 可以匹配全角标点
tom 和 jerry	\b 和 \b	否	\w 可以匹配中文字符
汤姆,jerry（半角标点）	\b 汤姆\b	能	\w 不能匹配半角标点
汤姆, jerry（全角标点）	\b 汤姆\b	否	\w 可以匹配全角标点

总的来说，如果使用 Unicode 匹配规则，尽量不要在处理中英文混排文本时使用 `\b`。如果使用 ASCII 匹配规则，则可以在处理英文文本时放心地使用 `\b`。

也有更复杂的情况，比如 Java 就是如此。在 Java 中，虽然 `\w` 只能匹配 `[0-9a-zA-Z_]`，`\b` 对“单词字符”的判断却是按照 Unicode 匹配规则的。`\b` 的匹配问题，本书后面关于各种语言的章节中都会专门讲解。

7.5 码值

Unicode 字符多种多样，除去 ASCII 中的字母、数字、标点和中文字符，还包括其他多种语言和多种符号，有些符号甚至很难打出来（比如表示商标注册的™），这时候该如何表示呢？再说远一点，如果我们想用一字符组匹配所有的中文字符，能不能像 `[a-z]` 那样呢？

使用正则表达式解决这类问题，必须依赖码值。前面讲过，每一个 Unicode 字符都有一个 Unicode 码值，所以在正则表达式中的 Unicode 字符往往采用 Unicode 码值来指定。

一般来说，指定码值的形式有两种：`\uxxxx` 和 `\u{xxxx}`（其中的 `xxxx` 为编码的值，`\u` 之后必须有 4 位十六进制数字）。.NET、Java、JavaScript 和 Python 使用前一种形式，而 PHP 和 Ruby 使用后一种形式（Ruby 1.9 以上版本才支持这种表示法，☞231；PHP 使用的字母是 `x` 而不是 `u`：`\x{xxxx}`）。

比如“发表”的“发”字对应的 Unicode 编码是 53 d1，所以可以在 .NET、Java、JavaScript 的正则表达式中用 `\u53d1` 表示它；Python 稍有不同，必须使用 `u"\u53d1"`（之前的 `u` 表示这是一个 Unicode 字符串，☞213）；在 Ruby 中必须写作 `\u{53d1}`；在 PHP 中则写作 `\x{53d1}`。表 7-5 总结了常用语言中用码值表示 Unicode 字符的方法。

表 7-5 常用语言中用码值表示 Unicode 字符的方法

编码	语言	表示法	说明
53 d1	.NET	\u53d1	
	Java	\u53d1	
	JavaScript	\u53d1	

( 续表 )

编码	语言	表示法	说明
53 d1	Python	\u53d1	必须使用 Unicode 字符串，在 Python 2 中，要在字符串之前加 u 指定 Unicode 字符串
	PHP	\x{53d1}	必须指定 Unicode 模式
	Ruby	\u{53d1}	限 Ruby 1.9 以上版本，且必须显式指定 Unicode 模式

既然可以这样指定 Unicode 字符，自然也可以在字符组中用范围表示法（☞105）指定一个 Unicode 编码范围了。这个功能，最常见的应用就是匹配任意一个中文字符。查询 Unicode 编码表可知，中文字符的编码点大多位于 4e00 到 9fff 之间<sup>1</sup>，所以可以用这样的字符组匹配中文字符（Unicode 编码 4e00~9fff 归类为 **CJK 统一表意符号 CJK Unified Ideographs**<sup>2</sup>，涵盖了绝大多数中文字符）。表 7-6 列出了常用语言中匹配中文字符的字符组。

表 7-6 常用语言中匹配中文字符的字符组

语言	字符组	说明
.NET	<code>[\\u4e00-\\u9fff]</code>	
Java	<code>[\\u4e00-\\u9fff]</code>	
JavaScript	<code>[\\u4e00-\\u9fff]</code>	
Python	<code>[\\u4e00-\\u9fff]</code>	必须使用 Unicode 字符串，在 Python 2 中，要在字符串之前加 u 指定 Unicode 字符串
PHP	<code>[\\x{4e00}-\\x{9fff}]</code>	必须指定 Unicode 模式
Ruby	<code>[\\u{4e00}-\\u{9fff}]</code>	限 Ruby 1.9 以上版本

如果实在没有办法使用 Unicode 编码环境，而只能采用 GBK 编码，在 Python、Ruby 1.8 和 PHP 中，也有个办法匹配所有中文字符：查阅 GBK 编码表可知，中文字符的 GBK 码值从 b0 00 开始，到 fe a0 结束，但是在通常情况下正则引擎无法正确识别非 Unicode 编码字符串的字符边界，所以不能写作 `[\\xb000-\\xfea0]`，只能写作 `[\\xb0-\\xfe][\\x00-\\xff]`，它匹配两个字节，第一个字节的码值范围从 b0 到 fe，第二个字节的码值范围从 00 到 ff，如果要匹配多个中文字符，必须添加括号将这两个字节分为一组，再使用量词，比如 `([\\xb0-\\xfe][\\x00-\\xff])`。

<sup>1</sup> 网上有许多资料称中文字符所处的字符范围在 4e00~9fa5 之间，而不是 4e00~9fff，这种说法有一定的道理，因为在 1992 年提交给 IRG（International Rapporteur Group）的字符只排到 9fa5。然而在这之后，制定 Unicode 4.2、5.1 和 5.2 规范时都进行了扩展，新增了字符。但是，4e00~9fff 是预留给东亚文字的编码点，所以使用 4e00~9fff 是更好的选择。具体信息，可以参考 <http://www.unicode.org/versions/Unicode6.0.0/ch12.pdf>。

<sup>2</sup> 请参考 [http://en.wikipedia.org/wiki/CJK\\_Unified\\_Ideographs](http://en.wikipedia.org/wiki/CJK_Unified_Ideographs)。

`xff]])+`。因为此时完全是将字符串作为单字节序列来对待的，所以也不应该指定任何与 Unicode 有关的设置（无论是正则表达式还是字符串，都是如此）。

## 7.6 Unicode 属性

每一个 Unicode 字符，除了有 Code Point 与之对应外，还具有其他属性，在正则表达式中常用到三种 Unicode 属性：Unicode Property、Unicode Block、Unicode Script，分别对应字符的功能、所属代码区段、书写系统；它们的表现形式都类似 `\p{property}`。图 7-3 比较形象地说明了这三种属性，下面详细介绍。

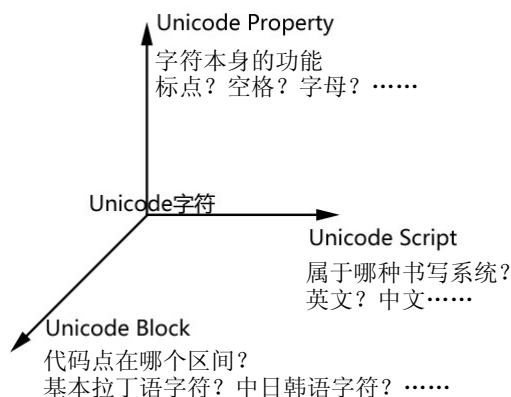


图 7-3 Unicode 字符的三种属性

### 7.6.1 Unicode Property

Unicode Property 的记法类似 `\p{L}`、`\p{P}`。它按照字符的功能分类 Unicode 字符，每个 Unicode 字符只能属于一个 Unicode Property。

可以这样理解 Unicode Property：它并不关心字符所属的语言，只关心字符的功能，比如 `\p{Z}` 表示任意的空白字符或不可见的分隔符；`\p{P}` 表示任何标点字符，等等。遇到中英文混排、全角、半角字符同时出现的情况，就可以用 `\p{Z}` 匹配所有的空白字符（而不用关心空格到底是全角空格还是半角空格），用 `\p{P}` 匹配所有的标点字符（而不用关心逗号到底是中文逗号还是英文逗号）。

如果把 Unicode Property 理解为一个“字符组”，那么一定还有对应的排除型字符组，此排除型字符组的通行记法是将 `\p{property}` 中的小写 `p` 改为大写 `P`，写作 `\P{property}`。这样，`\P{Z}` 对应 `\p{Z}` 无法匹配的字符，`\P{P}` 对应 `\p{P}` 无法匹配的字符。Unicode Block 和 Unicode Script 对应的排除型字符组也是这样标记，下面不再赘述。

支持 Unicode Property 的语言有 .NET、Java、PHP 和 Ruby (限 1.9 以上版本), 在 PHP 和 Ruby 中使用 Unicode Property 时, 必须要开启 Unicode 模式。下面的代码展示了在各种常用语言中, `\p{P}` 如何匹配一个全角逗号。

```
.NET
Regex.IsMatch(",", "\\p{P}"); // => True

Java
", ".matches("\\p{P}"); // => True

PHP
preg_match("/\\p{P}/u", ",", ""); // => 1

Ruby 1.9
/\\p{P}/u =~ ", " # => 0
```

Unicode Property 的列表附在章末。

## 7.6.2 Unicode Block

Unicode Block 不同于 Unicode Property, 它按照编码区间划分 Unicode 字符, 各个区间彼此联系但互不相交, 所以每个 Unicode 编码中的每个字符都有唯一归属的 Unicode Block。

在 Unicode 编码表中, 同一种语言的字符通常是落在同一区间的, 所以 Unicode Block 也可以粗略表示某类语言的字符, 比如 `\p{InHebrew}` 表示希伯来语字符, `\p{InCJK_Unified_Ideographs}` 表示兼容 CJK (中文、日文、韩文) 统一表意字符。如果细心阅读文档, 会发现 Unicode Block 的名字虽然类似某种语言的名字, 但都有 `In` (Java 风格) 或者 `Is` (.NET 风格) 前缀, 所以它对应的还是“落在某个区间的 Unicode 字符”。

在本书介绍的语言中, 只有 Java 和 .NET 支持 Unicode Block, 但是记法不相同, 下面以 CJK 统一表意字符为例, 虽然记法不同, 但功能是相同的, 都可以匹配任意一个中文字符。

```
Java
"我".matches("\\p{InCJK_Unified_Ideographs}"); // => True

.NET
Regex.IsMatch("我", "\\p{IsCJKUnifiedIdeographs}"); // => True
```

另外值得一用的是 `\p{InCJK_Symbols_and_Punctuation}`, 从命名上看, 它覆盖了中文、日文、韩文中的大部分常见标点, 但事实并非如此, 比如全角逗号, 就不在这个 Block 中。不过, 所有标点都在 `\p{P}` 这个 Unicode Property 中, 代码见例 7-14。

例 7-14 全角逗号不属于 `\p{InCJK_Symbols_and_Punctuation}`

```
". ".matches("\\p{InCJK_Symbols_and_Punctuation}"); // => True
", ".matches("\\p{InCJK_Symbols_and_Punctuation}"); // => False
```

```
". ".matches("\\p{P}");      // => True
", ".matches("\\p{P}");      // => True
```

注：Unicode Block 的列表附在章末。

7.6.3 Unicode Script

Unicode Script 按照字符所属的书写系统来划分 Unicode 字符，比如 `\p{Greek}` 表示希腊语字符，`\p{Han}` 表示汉语（中文字符）。它的写法类似 Unicode Block，只是名字的头没有 Is 或者 In。

在本书介绍的语言中，只有 PHP 和 Ruby（限 1.9 以上版本）支持 Unicode Script，PHP 在使用 Unicode Script 时，必须开启 Unicode 模式（详见 PHP 的章节）。在这两种语言中，我们可以很方便地用 `\p{Han}` 来匹配中文字符。

```
PHP
preg_match("/\p{Han}/u", "我");      // => 1
Ruby 1.9
/\p{Han}/u =~ "我"      # => 0
```

注：Unicode Script 的列表附在章末。

7.7 Unicode 属性列表

7.7.1 Unicode Property

每个 Unicode 字符都只能属于一个 Unicode Property。所有的 Unicode Property 共分为 7 大类，30 小类。大类的名字只有 1 个字母，小类的名字则不止 1 个字母，开头字母与所在大类的名字相同，小类包含的字符都属于它所在的大类。表 7-7 列出了所有的 Unicode Property。

表 7-7 Unicode Property

Unicode Property		说明
\p{C}		不可见的控制字符和未使用的码值
	\p{Cc}	ASCII 编码中 0x00 到 0x1F 或 Latin-1 编码中 0x80 到 0x9F 的控制字符
	\p{Cf}	不可见的格式字符
	\p{Co}	留作私用的码值
	\p{Cs}	UTF-16 编码中 surrogate pair 的一半
	\p{Cn}	未指定的码值



(续表)

Unicode Property		说明
\p{L}		各种语言中的字母
	\p{Ll}	具有大写形式的字母的小写形式
	\p{Lt}	只有在单词首位才大写的字符
	\p{L&}	等于 Ll、Lu、Lt 的组合
	\p{Lo}	没有大小写形态的字符
	\p{Lu}	具有小写形式的字母的大写形式
\p{M}		用来与其他字符结合的字符（声调、元音变化符等）
	\p{Mc}	与其他字符组合，并且会占用空间的字符（常见于东亚语言）
	\p{Me}	需要成对出现的字符，比如圆括号、方括号
	\p{Mn}	用来与其他字符结合，但并不占用额外空间的字符
\p{N}		各种书写系统中的数字字符
	\p{Nd}	各种书写系统中的 0~9 的字符
	\p{Nl}	形如字符的数字，比如罗马数字
	\p{No}	上标或者下标数字，或者是 0~9 之外的数字（不包括表意书写系统中的数字）
\p{P}		各种标点符号
	\p{Pd}	各种连字符号
	\p{Ps}	成对但不同的符号的前半部分（包括英文括号、中文括号、书名号）
	\p{Pe}	成对但不同的符号的后半部分（包括英文括号、中文括号、书名号）
	\p{Pi}	成对且相同的符号的前半部分（比如单引号、双引号）
	\p{Pf}	成对且相同的符号的后半部分（比如单引号、双引号）
	\p{Pc}	类似下画线之类的标点字符
	\p{Po}	除横线、括号、引号和连接符之外的任何标点符号
\p{S}		数字符号、货币符号
	\p{Sm}	数字符号
	\p{Sc}	货币符号
	\p{Sk}	由多个字符构成的组合字符
	\p{So}	数字符号、货币符号和组合字符之外的符号字符
\p{Z}		空白字符，或者不可见的分隔符
	\p{Zs}	不可见但占用空间的空白字符
	\p{Zl}	分行符 U+2028
	\p{Zp}	分段符 U+2029

## 7.7.2 Unicode Block

每个 Unicode Block 都对应一个连续的 Unicode 码值区间，U+0000 到 U+FFFF 之间的字符一共划分为 105 个 Block。

使用时应该注意，Java 使用的 Unicode Block 是 `\p{In...}` 形式的<sup>1</sup>，比如 `InCJKUnifiedIdeographs`；而 .NET 使用的 Unicode Block 是 `\p{Is...}` 形式的，同时不会包含下画线<sup>2</sup>，比如 `IsCJKUnifiedIdeographs`。表 7-8 列出了所有的 Unicode Block，以 Java 风格为例<sup>3</sup>。

表 7-8 Unicode Block

码值	Unicode Block	说明
U+0000~U+007F	<code>\p{InBasic_Latin}</code>	控制符及基本拉丁文
U+0080~U+00FF	<code>\p{InLatin_1_Supplement}</code>	控制符及拉丁文补充
U+0100~U+017F	<code>\p{InLatin_Extended-A}</code>	拉丁文扩展-A
U+0180~U+024F	<code>\p{InLatin_Extended-B}</code>	拉丁文扩展-B
U+0250~U+02AF	<code>\p{InIPA_Extensions}</code>	国际音标扩展
U+02B0~U+02FF	<code>\p{InSpacing_Modifier_Letters}</code>	空白修饰字符
U+0300~U+036F	<code>\p{InCombining_Diacritical_Marks}</code>	结合用读音符号
U+0370~U+03FF	<code>\p{InGreek_and_Coptic}</code>	希腊文及科普特文
U+0400~U+04FF	<code>\p{InCyrillic}</code>	西里尔字母
U+0500~U+052F	<code>\p{InCyrillic_Supplementary}</code>	西里尔字母补充
U+0530~U+058F	<code>\p{InArmenian}</code>	亚美尼亚语
U+0590~U+05FF	<code>\p{InHebrew}</code>	希伯莱语
U+0600~U+06FF	<code>\p{InArabic}</code>	阿拉伯语
U+0700~U+074F	<code>\p{InSyriac}</code>	叙利亚语
U+0780~U+07BF	<code>\p{InThaana}</code>	马尔代夫语
U+0900~U+097F	<code>\p{InDevanagari}</code>	天城文书
U+0980~U+09FF	<code>\p{InBengali}</code>	孟加拉语
U+0A00~U+0A7F	<code>\p{InGurmukhi}</code>	锡克教文
U+0A80~U+0AFF	<code>\p{InGujarati}</code>	古吉拉特文
U+0B00~U+0B7F	<code>\p{InOriya}</code>	奥里亚文

<sup>1</sup> 具体情况可参考 <http://download.oracle.com/javase/6/docs/api/java/lang/Character.UnicodeBlock.html>。

<sup>2</sup> 具体情况可参考 <http://msdn.microsoft.com/en-us/library/20bw873z.aspx#SupportedNamedBlocks>。

<sup>3</sup> 如果想查看这些 Unicode Block 的具体内容，可以参考 <http://www.nengcha.com/code/unicode/class/>。

(续表)

码值	Unicode Block	说明
U+0B80~U+0BFF	\p{InTamil}	泰米尔文
U+0C00~U+0C7F	\p{InTelugu}	泰卢固文
U+0C80~U+0CFF	\p{InKannada}	卡纳达文
U+0D00~U+0D7F	\p{InMalayalam}	德拉维族语
U+0D80~U+0DFF	\p{InSinhala}	僧伽罗语
U+0E00~U+0E7F	\p{InThai}	泰语
U+0E80~U+0EFF	\p{InLao}	老挝语
U+0F00~U+0FFF	\p{InTibetan}	藏语
U+1000~U+109F	\p{InMyanmar}	缅甸语
U+10A0~U+10FF	\p{InGeorgian}	格鲁吉亚语
U+1100~U+11FF	\p{InHangul_Jamo}	朝鲜语
U+1200~U+137F	\p{InEthiopic}	埃塞俄比亚语
U+13A0~U+13FF	\p{InCherokee}	切诺基语
U+1400~U+167F	\p{InUnified_Canadian_Aboriginal_Syllabics}	统一加拿大土著语音节
U+1680~U+169F	\p{InOgham}	欧甘字母
U+16A0~U+16FF	\p{InRunic}	如尼文
U+1700~U+171F	\p{InTagalog}	塔加拉语
U+1720~U+173F	\p{InHanunoo}	Hanunóo
U+1740~U+175F	\p{InBuhid}	Buhid
U+1760~U+177F	\p{InTagbanwa}	Tagbanwa
U+1780~U+17FF	\p{InKhmer}	高棉语
U+1800~U+18AF	\p{InMongolian}	蒙古语
U+1900~U+194F	\p{InLimbu}	Limbu
U+1950~U+197F	\p{InTai_Le}	德宏泰语
U+19E0~U+19FF	\p{InKhmer_Symbols}	高棉语记号
U+1D00~U+1D7F	\p{InPhonetic_Extensions}	语音学扩展
U+1E00~U+1EFF	\p{InLatin_Extended_Additional}	拉丁文扩充附加
U+1F00~U+1FFF	\p{InGreek_Extended}	希腊语扩展
U+2000~U+206F	\p{InGeneral_Punctuation}	常用标点
U+2070~U+209F	\p{InSuperscripts_and_Subscripts}	上标及下标
U+20A0~U+20CF	\p{InCurrency_Symbols}	货币符号

(续表)

码值	Unicode Block	说明
U+20D0~U+20FF	\p{InCombining_Diacritical_Marks_for_Symbols}	组合用记号
U+2100~U+214F	\p{InLetterlike_Symbols}	字母式符号
U+2150~U+218F	\p{InNumber_Forms}	数字形式
U+2190~U+21FF	\p{InArrows}	箭头
U+2200~U+22FF	\p{InMathematical_Operators}	数学运算符
U+2300~U+23FF	\p{InMiscellaneous_Technical}	杂项工业符号
U+2400~U+243F	\p{InControl_Pictures}	控制图片
U+2440~U+245F	\p{InOptical_Character_Recognition}	光学识别符
U+2460~U+24FF	\p{InEnclosed_Alphanumerics}	封闭式字母和数字
U+2500~U+257F	\p{InBox_Drawing}	制表符号
U+2580~U+259F	\p{InBlock_Elements}	方块元素
U+25A0~U+25FF	\p{InGeometric_Shapes}	几何图形
U+2600~U+26FF	\p{InMiscellaneous_Symbols}	杂项符号
U+2700~U+27BF	\p{InDingbats}	印刷符号
U+27C0~U+27EF	\p{InMiscellaneous_Mathematical_Symbols-A}	杂项数学符号 A
U+27F0~U+27FF	\p{InSupplemental_Arrows-A}	箭头补充 A
U+2800~U+28FF	\p{InBraille_Patterns}	盲点文字模型
U+2900~U+297F	\p{InSupplemental_Arrows-B}	箭头补充 B
U+2980~U+29FF	\p{InMiscellaneous_Mathematical_Symbols-B}	杂项数学符号 B
U+2A00~U+2AFF	\p{InSupplemental_Mathematical_Operators}	数学运算符补充
U+2B00~U+2BFF	\p{InMiscellaneous_Symbols_and_Arrows}	杂项箭头和符号
U+2E80~U+2EFF	\p{InCJK_Radicals_Supplement}	CJK 部首补充
U+2F00~U+2FDF	\p{InKangxi_Radicals}	康熙字典部首
U+2FF0~U+2FFF	\p{InIdeographic_Description_Characters}	表意文字描述符
U+3000~U+303F	\p{InCJK_Symbols_and_Punctuation}	CJK 标点和符号
U+3040~U+309F	\p{InHiragana}	日文平假名
U+30A0~U+30FF	\p{InKatakana}	日文片假名
U+3100~U+312F	\p{InBopomofo}	注音字母
U+3130~U+318F	\p{InHangul_Compatibility_Jamo}	朝鲜文兼容字母
U+3190~U+319F	\p{InKanbun}	象形字注释标注
U+31A0~U+31BF	\p{InBopomofo_Extended}	注音字母扩展
U+31F0~U+31FF	\p{InKatakana_Phonetic_Extensions}	日文片假名语音扩展

(续表)

码值	Unicode Block	说明
U+3200~U+32FF	\p{InEnclosed_CJK_Letters_and_Months}	封闭式 CJK 文字和月份
U+3300~U+33FF	\p{InCJK_Compatibility}	CJK 兼容
U+3400~U+4DBF	\p{InCJK_Unified_Ideographs_Extension_A}	CJK 统一表意符号扩展 A
U+4DC0~U+4DFF	\p{InYijing_Hexagram_Symbols}	易经六十四卦符号
U+4E00~U+9FFF	\p{InCJK_Unified_Ideographs}	CJK 统一表意符号
U+A000~U+A48F	\p{InYi_Syllables}	彝文音节
U+A490~U+A4CF	\p{InYi_Radicals}	彝文字根
U+AC00~U+D7AF	\p{InHangul_Syllables}	朝鲜文音节
U+D800~U+DB7F	\p{InHigh_Surrogates}	UTF-16 高半区域
U+DB80~U+DBFF	\p{InHigh_Private_Use_Surrogates}	自行使用的 UTF-16 高半区域
U+DC00~U+DFFF	\p{InLow_Surrogates}	UTF-16 低半区域
U+E000~U+F8FF	\p{InPrivate_Use_Area}	自行使用区域
U+F900~U+FAFF	\p{InCJK_Compatibility_Ideographs}	CJK 兼容象形文字
U+FB00~U+FB4F	\p{InAlphabetic_Presentation_Forms}	字母表达形式
U+FB50~U+FDFF	\p{InArabic_Presentation_Forms-A}	阿拉伯表达形式 A
U+FE00~U+FE0F	\p{InVariation_Selectors}	变量选择符
U+FE20~U+FE2F	\p{InCombining_Half_Marks}	组合用半符号
U+FE30~U+FE4F	\p{InCJK_Compatibility_Forms}	CJK 兼容形式
U+FE50~U+FE6F	\p{InSmall_Form_Variants}	小型变体形式
U+FE70~U+FEFF	\p{InArabic_Presentation_Forms-B}	阿拉伯表达形式 B
U+FF00~U+FFEF	\p{InHalfwidth_and_Fullwidth_Forms}	半型及全型形式
U+FFF0~U+FFFF	\p{InSpecials}	特殊

### 7.7.3 Unicode Script

除去没有赋值的 Unicode 码值之外, 每个 Unicode 码值都归属于某个 Unicode Script, 它一般对应某种语言。表 7-9 列出了所有的 Unicode Script。

表 7-9 Unicode Script

\p{Common}	\p{Georgian}	\p{Limbu}	\p{Han}
\p{Arabic}	\p{Inherited}	\p{Malayalam}	\p{Tamil}
\p{Armenian}	\p{Greek}	\p{Mongolian}	\p{Telugu}
\p{Bengali}	\p{Kannada}	\p{Myanmar}	\p{Thaana}
\p{Bopomofo}	\p{Gujarati}	\p{Ogham}	\p{Thai}

(续表)

\p{Braille}	\p{Katakana}	\p{Oriya}	\p{Tibetan}
\p{Buhid}	\p{Gurmukhi}	\p{Runic}	\p{Yi}
\p{CanadianAboriginal}	\p{Khmer}	\p{Sinhala}	
\p{Cherokee}	\p{Han}	\p{Syriac}	
\p{Cyrillic}	\p{Lao}	\p{Tagalog}	
\p{Devanagari}	\p{Hangul}	\p{Tagbanwa}	
\p{Ethiopic}	\p{Latin}	\p{TaiLe}	

## 7.8 POSIX 字符组

在第 1 章, 我们已经介绍过 POSIX 字符组, 并且提到, 根据 locale 的不同, 这些字符组简记法能匹配的字符也不相同, 表 7-10 介绍了 Unicode 编码中这些简记法能匹配的字符。

表 7-10 Unicode 编码中的 POSIX 字符组

POSIX	说明	Unicode 属性表示法
<code>[ :alnum: ]</code>	字母字符和数字字符	<code>[ \p{L&amp;} \p{Nd} ]</code>
<code>[ :alpha: ]</code>	字母	<code>\p{L&amp;}</code>
<code>[ :ascii: ]</code>	ASCII 字符	<code>\p{InBasicLatin}</code>
<code>[ :blank: ]</code>	空格字符和制表符	<code>[ \p{Zs} \t ]</code>
<code>[ :cntrl: ]</code>	控制字符	<code>\p{Cc}</code>
<code>[ :digit: ]</code>	数字字符	<code>\p{Nd}</code>
<code>[ :graph: ]</code>	空白字符之外的字符	<code>[ ^\p{Z} \p{C} ]</code>
<code>[ :lower: ]</code>	小写字母字符	<code>\p{Ll}</code>
<code>[ :print: ]</code>	类似 <code>[ :graph: ]</code> , 但包括空白字符	<code>\p{C}</code>
<code>[ :punct: ]</code>	标点符号	<code>[ \p{P} \p{S} ]</code>
<code>[ :space: ]</code>	空白字符	<code>[ \p{Z} \t \r \n \v \f ]</code>
<code>[ :upper: ]</code>	大写字母字符	<code>\p{Lu}</code>
<code>[ :word: ]*</code>	字母字符	<code>[ \p{L} \p{N} \p{Pc} ]</code>
<code>[ :xdigit: ]</code>	十六进制字符, 等于 <code>[ 0-9a-zA-Z ]</code>	<code>[ A-Fa-f0-9 ]</code>

注: 标记\*的字符组简记法并不是 POSIX 规范中的, 但使用很多, 一般语言中都提供了, 文档也会提到。

在本书介绍的 6 种语言中, Java、PHP、Ruby 支持使用 POSIX 字符组。但是只有 Ruby 1.9 的 POSIX 字符组是完全支持 Unicode 字符的, 并且不需要显式指定 Unicode 模式。在 Java 和 PHP 中, POSIX 字符组只能匹配 ASCII 字符。

## 第8章 匹配原理

在固定字符串的处理上，正则表达式的速度是赶不上简单字符串处理的；如果要进行复杂多变的字符串处理，正则表达式的速度则要胜于简单的字符串处理，比如正则表达式 `a(bb)+a`，它能匹配的字符串是 `abba`、`abbbba`、`abbbbbba`...，如果用简单的字符串匹配算法，我们可能需要逐个寻找可能的形态，也就是 `abba`、`abbbba`、`abbbbbba`...，此类尝试不知道要进行多少次，效率自然大打折扣。这难免让人好奇，正则表达式为什么这样快呢？要弄清这个问题，必须详细了解正则表达式的匹配原理<sup>1</sup>。

### 8.1 有穷自动机

正则表达式能迅速进行复杂处理的秘密就在于，它采用了一种特殊的理论模型：**有穷自动机**（Finite Automata，也叫**有穷状态自动机**，finite-state machine）。这种机器具备有限个状态，可以根据不同的条件在状态之间转移。

卖饮料的自动售货机就是一种有穷自动机：假设其中的饮料价格都是整数元，而且只接收5块钱的纸币，根据余额的不同，可能状态有6个：5元、4元、3元、2元、1元、0元。你塞进去5块钱，此时的状态就是“5元”，你点了一罐可乐，花去3元，于是状态切换到“2元”，这时候你按下“退币”，就把剩下的2元退给你，并把状态切换到“0元”，“0元”这个状态也叫做“最终状态”。到此，这一轮状态转移结束，如果你再塞5块钱，就开始新一轮的状态转移。

严格说起来，有穷自动机必须满足4个条件：

- （1）具有有限多个状态；
- （2）有一套状态转移函数（或者叫“规则”）；
- （3）有一个开始状态；
- （4）有一个或多个最终状态。

我们说自动售货机是一种有穷自动机，就是因为它满足这4个条件：

---

<sup>1</sup> 本章关于自动机原理的部分示意图来自在 Google 工作的 Russ Cox (<http://swtch.com/~rsc/>)，已经获得作者本人许可。

- (1) 具有有限多个状态（6个）；
- (2) 有一套状态转移函数（比如余额还有3元，你买了一罐2元的饮料，则转移到状态“1元”，如果你选择买4元的饮料，则报告“余额不足”，状态并不变化）；
- (3) 有一个开始状态（余额“5元”）；
- (4) 有一个最终状态（余额“0元”）。

自动售货机对应的有穷自动机模型如图 8-1 所示，它包含 6 个状态，对应余额的 6 种可能，起始状态是“¥5”，结束状态是“¥0”，每个箭头代表一个转移函数（每样商品的价格为 1 元或者 2 元，所以每个转移函数上的文字或者是-1，或者是-2），在¥4、¥3、¥2、¥1 状态下，都可以直接退币，所以有一个转移函数直达最终状态。

自动售货机并不关心你买了什么商品，也不关心你的选择顺序，无论你买什么，它总处在这 6 个状态之一，只需要根据状态转移函数在其中转移即可。

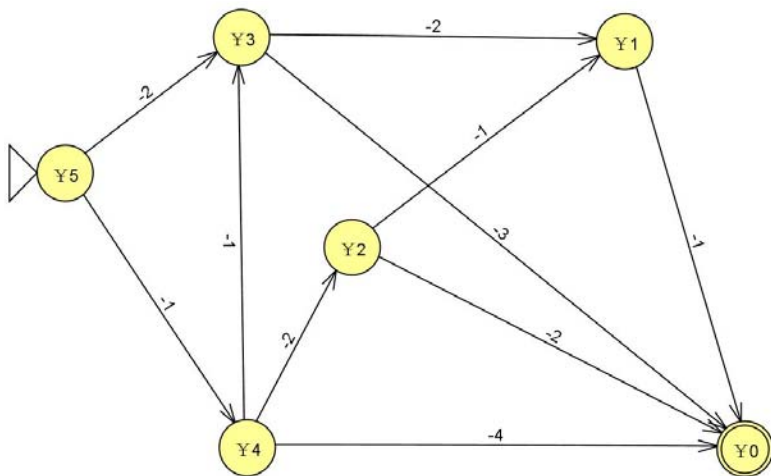


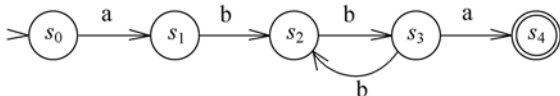
图 8-1 自动售货机对应的有穷自动机模型

## 8.2 正则表达式的匹配过程

正则表达式所使用的理论模型就是有穷自动机，其具体实现称为**正则引擎**（Regex Engine）。用正则表达式处理字符串，首先需要生成自动机（你应该还记得，很多语言中使用正则表达式之前都要“编译”正则对象）；之后，无论输入什么字符串，正则引擎都只需要老老实实地在状态之间游走。

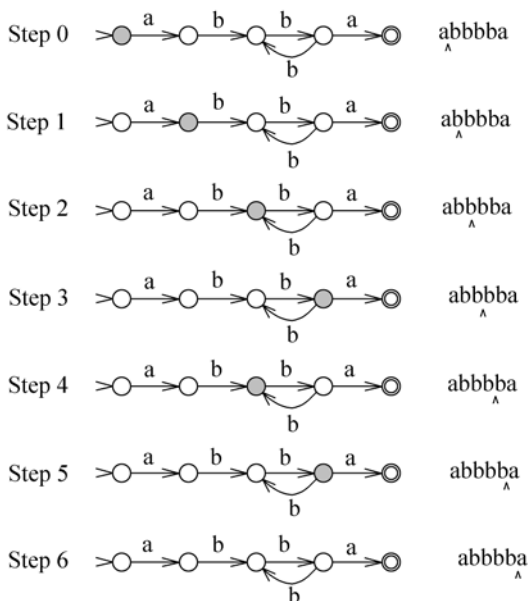
图 8-2 显示了正则表达式 `a(bb)+a` 对应的自动机。这台自动机的表示与之前看到的稍有不同：在匹配字符串时，输入的都是字符，所以箭头上标注的都是字符。



图 8-2  $a(bb)^+a$  对应的有穷自动机

在这台有穷自动机中,  $s_0$ 、 $s_1$ 、 $\dots$ 、 $s_4$  是各个状态,  $s_0$  为开始状态,  $s_4$  为最终状态; 转移函数很直观: 比如当前状态是  $s_0$ , 输入字符 **a**, 则转移到  $s_1$ ; 如果当前状态为  $s_0$ , 输入的不是 **a**, 那么直接退出。这也很好理解: 如果正则表达式是  $a(bb)^+a$ , 它能匹配的字符串只能是以字符 **a** 开头的, 否则必然不能匹配。

图 8-3 说明了这台有穷自动机对字符串 **abbbba** 的处理过程。

图 8-3 自动机对字符串 **abbbba** 的匹配过程

在经历了一系列的状态转移之后, 字符串 **abbbba** 处理完毕, 自动机停留在最终状态上, 也就是说, 字符串 **abbbba** 可以由正则表达式  $a(bb)^+a$  匹配。

同一个正则表达式对应有穷自动机不止一台, 可以是若干台, 这些有穷自动机是等价的。同样是正则表达式  $a(bb)^+a$ , 它对应到图 8-4 所示的两台完全等价的有穷自动机。

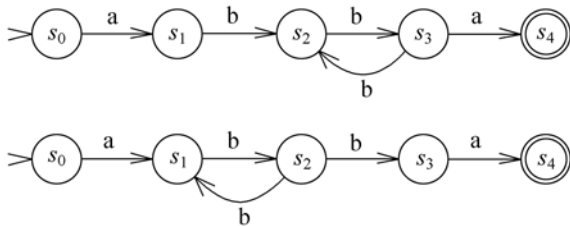


图 8-4 两台完全等价的自动机

仔细观察会发现，第二台自动机有些奇怪，在输入 **ab** 之后，再输入 **b**，它所处的状态是不确定的：可能在  $S_1$ ，也可能在  $S_3$ 。但是，输入  $a(bb)+a$  能匹配的字符串，它确实可以抵达最终状态  $S_4$ 。图 8-5 所示的自动机看起来更加奇怪，而且它仍然是与  $a(bb)+a$  完全对应的。

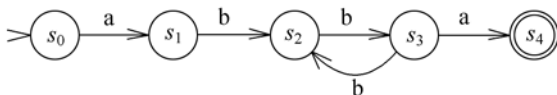


图 8-5 另一台等价的自动机

在状态  $S_3$ ，即便没有输入任何字符，也不会停留下来，而可能“凭空”转义到  $S_1$ 。也就是说，在某个时刻，自动机到底处在状态  $S_3$ ，还是  $S_1$ ，这是不确定的！但是这种不确定性，并不会影响自动机对于正则表达式  $a(bb)+a$  的匹配。也就是说，这台有穷自动机与之前的两台有穷自动机，也是完全等价的！

根据状态的确定与否，一般我们会把有穷自动机（正则引擎）分为两类：一类是**确定型有穷自动机**（Definite Finite Automata，简称 DFA），在任何时刻，它所处的状态是确定无疑的；另一类是**非确定型有穷自动机**（Non-definite Finite Automata，简称 NFA），在某个时刻，它所处的状态可能是不确定的。图 8-6 把上面的三台自动机做了分类，第一台是 DFA，而另两台是 NFA。

DFA:



NFA:

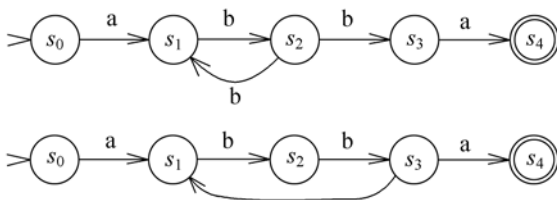


图 8-6 等价的 DFA 和 NFA

可以证明，DFA 和 NFA 之间存在等价关系；也就是说，每一台 DFA 都可以等价转换为一台 NFA，反过来也成立（具体证明过程很复杂，这里略去，有兴趣的读者可以参考计算理论相关资料）。

比较正则表达式  $a(bb)+a$  和这三台自动机，会发现 NFA 构造起来更直观，实际上这是普遍规律：从正则表达式出发，构造 NFA 的难度要小于 DFA。但是正如之前讲过的，DFA 在任意时刻必定处于某个确定的状态，而 NFA 可能处于若干状态之中的任何一个，所以，如果使用 NFA，就必须保存所有的可能的状态，并且在某种状态不可行时“回退”到之前保存的状态，这就是正则表达式匹配中的重要概念：回溯。

## 8.3 回溯

比起 DFA，NFA 看起来足够“麻烦”：它的状态是不确定的，这有点像走迷宫，越走岔路口越多，最后不会迷路吗？

不过，NFA 的正则引擎自有办法：如果有多个可能的状态，它们会在选择时记录下这些状态备用，然后才选择其中某个状态尝试；如果之后遇到死路，则退回去，选择最近一次记录的且未尝试过的状态；如果又遇到死路，再选择最近一次记录的且未尝试过的状态……这有点像在岔路口留下标记——如果我们在遇到的每个岔路口都留下标记，即便前头是死路，也可以根据标记返回，而不会迷路。

为了说明 NFA 的匹配过程，来看在第 2 章举过的双引号字符串匹配的例子，所用的正则表达式是 `".*"`，而字符串是 `"quoted string"`，匹配的过程如图 8-7 所示。

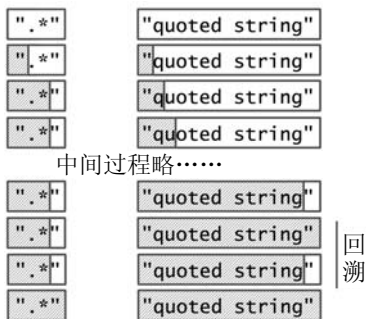


图 8-7 表达式 `".*"` 对字符串 `"quoted string"` 的匹配过程

从图中可以看出，在匹配的过程中，`.*` 曾经匹配了 `quoted string`，但为了保证表达式中最后一个 `"` 的匹配，`.*` 不得不“交还”最后的 `"`，这种“尝试失败-重新选择”的过程，就是回溯（backtracking）。

回溯只属于 NFA 引擎。从之前的原理图中可以看到，NFA 匹配时，正则引擎并不准确知道当前的状态，只能在所有状态不确定的地方将各种状态都保存下来（现在已经匹配了哪些字符，进行到字符串中的哪个位置，正则表达式中的哪个位置），逐一尝试，发现此路不通，则退回来，选择最近保存的其他状态尝试……如此持续进行下去，直到达到最终状态（这时候报告“在整个正则表达式开始尝试的位置，匹配成功”）；或者所有可能状态都尝试完毕，仍然不能到达最终状态（如果当前位置是字符串的末尾，则报告“在当前位置匹配失败”；否则，把“整个正则表达式开始尝试的位置”向前推进一个字符，再开始新一轮的尝试）。

看到这里，就不难明白为什么不推荐使用 `.*` 了，因为 `.` 几乎能匹配任何字符串（如果明确指定单行模式，则确实能匹配任何字符），而 `*` 又表示“匹配优先”，所以正则引擎在处理 `.*` 后的其他元素之前，会先让 `.*` “吞掉”几乎整个字符串。仍然是上面的正则表达式，只是字符串变为

"quoted" string，回溯的次数大大增加了，如果在结尾的"之后还有很长的文本，回溯的次数还可能大大增加，匹配过程如图 8-8 所示。

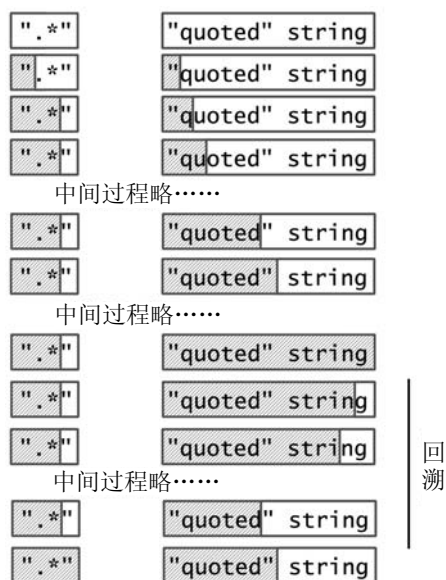


图 8-8 表达式 `.*` 对字符串 `"quoted" string` 的匹配过程

为避免这类问题，最好的办法是准确表达意图，比如规定双引号字符串内部不容许出现双引号字符，就要将表达式改为 `"[^"]*"`；当然也可以换用忽略优先量词，将表达式改为 `".*?"`，两种办法都可行<sup>1</sup>。不过总的经验是，除非确实必要，否则尽量不要使用 `.*`。

要注意的不仅仅有 `.*`，还有更遭的情况，比如 `(...)*` 之类的表达式，这时候回溯的次数会呈指数增长，却不会对匹配有任何影响，所以应该绝对避免。第 3 章中匹配 HTML tag 的正则表达式是 `<('[^']*'|\"[^\"]*"|'>)+>`，其中的多选分支 `[^'\">]` 没有添加量词 `*`，就是因为单引号字符串和双引号字符串之外的字符虽然可能有很多，但多选结构最外层还有 `+` 限定，从忽略之前两个多选分支来看，`([^'\">]+)` 要好过 `([^'\">]*)`。关于这类问题的详细讲解，可以参考《精通正则表达式》（第 3 版）的第 4 章相关内容。

在实际应用中，不只要注意自己写的正则表达式，还需要防范外界的恶意程序，它们刻意使用会造成大量回溯的表达式，将计算机的资源消耗殆尽，这种攻击有一个专门的名词，叫做正则表达式拒绝服务攻击（Regular Expression Denial of Service）<sup>2</sup>。

<sup>1</sup> 其实这个表达式并不够完整，因为有些双引号字符串中是可以出现双引号字符的，比如 `"\""`，其中的双引号字符之前有反斜线，所以它是一个经过转义的双引号。能处理这种情况的表达式在**错误！未定义书签**。页详细讨论。

<sup>2</sup> 详细信息可以参考 <http://en.wikipedia.org/wiki/ReDoS>。

## 8.4 NFA 和 DFA

上一节粗略介绍了回溯，它是 NFA 特有的功能，DFA 不需要回溯，也就不需要保存状态，再反复尝试。这样看来，NFA 不是要更慢吗？事实也确实如此，但是当前我们所使用的大多数工具中的正则引擎，都选用了 NFA，这是为什么呢？

NFA 确实更慢，但 NFA 也有自己的优势：如果正则表达式比较复杂，构建 NFA 的时间比 DFA 的时间更短（举例来说，如果你的正则表达式使用了多选分支，每个分支其实只是个简单的字符串，那么你完全可以直接对每个多选分支构建简单的 NFA，再把它们简单“并列”起来就可以了；相比之下，构建整个表达式对应的 DFA 就复杂多了）。

同时，现代 NFA 也提供了更多的优化措施，比如之前提到的 `a(bb)+a` 的匹配，优化过的 NFA 可以“并行尝试”，其匹配过程如图 8-9 所示，这样的速度就快多了。

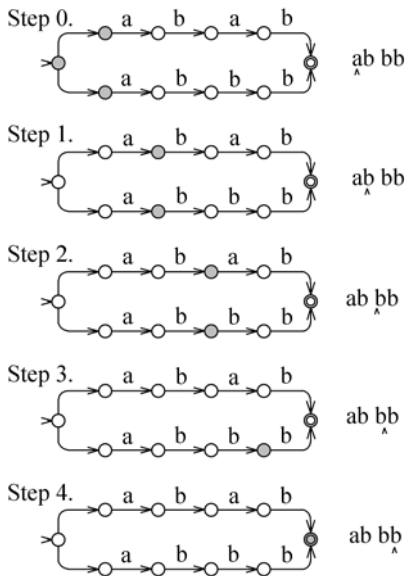


图 8-9 经过优化的 NFA 匹配过程

更重要的是，NFA 的匹配性质决定了它必须在匹配过程中保存可能的状态，需要“停下来四处看看”，所以也能够“回顾一路走来的历程”；相比之下，DFA 不会两次测试同一个字符，所以不需要保存状态。因此，NFA 具有许多 DFA 无法提供的功能：比如捕获型括号 `(...)`，反向引用 `\num`，环视功能 `(?!...)`、`(?=...)`，忽略优先量词 `+?`、`*?`、`??`……

如果希望用到这些功能，一定不要选择使用 DFA 引擎的工具。当然一般来说，用户并不需要操心引擎是 DFA 或者 NFA，毕竟它们是位于“幕后”的，需要关注只是，是否提供了希望实

现功能所用到的 API。

而且，现代的一些工具中，为兼顾效率和功能，同时包含了 DFA 和 NFA 两种引擎，如果发现正则表达式中没有专属于 NFA 的功能，则使用 DFA，否则使用 NFA。表 8-1 列出了各种常用工具所使用的正则引擎。

表 8-1 常用工具所使用的正则引擎

引擎类型	程序
DFA	awk（大多数版本）、egrep（大多数版本）、flex、lex、MySQL、Procmail
NFA	GNU Emacs、Java、grep、less、more、.NET、Python、Ruby、PHP、sed

细分起来，NFA 又有传统型 NFA（Traditional NFA）和 POSIX NFA 两种。两者的主要区别在于，如果多选分支中的多个分支都能匹配，传统型 NFA 优先选择左侧的分支，而 POSIX NFA 一定要选择最长的分支。比如用表达式 `(jeff|jeffrey)` 匹配字符串 `jeffrey`，POSIX NFA 的结果是 `jeffrey`，传统型 NFA 的结果则是 `jeff`——如果调换多选分支的顺序，写成 `(jeffrey|jeff)`，POSIX NFA 的结果不变，传统型 NFA 的结果则变为 `jeffrey`。

问题看起来很复杂，具体使用起来其实比较简单：POSIX NFA 的应用很少，主要用于 Linux/UNIX 下的工具（所以它们中的很多并不支持捕获分组，详见第 16 章），本书中介绍的编程语言基本都采用传统型 NFA 引擎。保险起见，不妨这样记忆：一般情况下，多选分支优先采用最左侧的分支。

关于正则表达式的匹配原理就介绍这么多，如果对这部分内容有更深入的兴趣，可以参考《精通正则表达式》（第 3 版）的第 4 章。

# 第 9 章 常见问题的解决思路

前面已经讲解了正则表达式的常用功能和结构，Unicode 问题，正则表达式的匹配原理，关于正则表达式的通用知识已经基本介绍完毕。不过，使用正则表达式是为了解决各种具体问题，所以这一章专门讲解使用正则表达式解决具体问题的思路和注意事项。

## 9.1 关于元素的三种逻辑

正则表达式提供了一整套描述文本特征的方法，它的匹配其实也就是查找符合所描述特征的文本的过程。按照元素（单个字符、字符组、多选分支等等）的出现情况，这些特征分为三类，也可称为三种逻辑：必须出现、可能出现、不能出现；具体解释如表 9-1。

表 9-1 关于元素的三种逻辑

分类	说明
必须出现	某些元素必须出现
可能出现	某个元素或许出现，或许不出现，或许长度不固定；要出现的，是某几个元素中的一个
不能出现	某些元素不能出现

不管正则表达式多么复杂，总是这三种逻辑的组合。比如匹配双引号字符串的任务，可以按三种逻辑分析如下。

必须出现	首尾的双引号字符必须出现
可能出现	两个双引号之间的字符个数是不确定的（如果是空字符串""，则两个双引号之间没有字符）
不能出现	两个双引号之间不容许出现双引号字符

再比如，匹配 HTML 中的 open tag（如<h1>）和 close tag（如</h1>），按三种逻辑分析如下。

必须出现	首尾必须分别是<和>，如果是 close tag，则<-之后必须出现/
可能出现	<和>-之间字符串长度必须大于 1（<>-不是一个合法的 tag）
不能出现	<和>-之间不能有>，如果是 open tag，<-之后不能出现/

9.1.1 必须出现

“必须出现”是正则表达式中最普通的逻辑关系，它表示某个元素必须出现。通常，这些元素是所要匹配文本的最重要特征：查找 `tag`，则必须出现的是字符`<`和`>`；查找 `E-mail` 地址，必须出现的元素是`@`。

如果某个元素必须出现，通常不会（也不应该）用量词（`*`、`?`）来限制，也不出现在多选结构中（如果把普通的字符串也看做正则表达式，那么其中每个字符都是必须出现的）。所以，如果在匹配 `tag` 的正则表达式中，`<`和`>`出现在某个多选结构内，或者之后跟有`?`、`*`之类的量词，那么这个表达式多半有问题；同样，在匹配 `E-mail` 地址的正则表达式中的`@`字符，也不应该有量词限定，或者出现在多选结构内部。

这条要求看起来简单，但是在日常使用中却容易犯错误。经常遇到的情况是，为了考虑更多的可能情况而修改正则表达式，增加量词和多选结构，最后正则表达式中没有任何必须出现的元素。比如匹配数字字符串的正则表达式，一个数字字符串可能包含三个部分：开头的`+`或`-`、整数部分、小数点和小数部分，然后分别列出匹配这三个部分的正则表达式。

符号	<code>[-+]</code>
整数部分	<code>\d+</code>
小数部分	<code>\.\d+</code>

单独来看，这三个部分都不是必须出现的：数字字符串可以没有符号，比如 `3.14`；也可以没有整数部分，比如`.14`；还可以没有小数点和小数部分，比如`-3`。所以，有些人就把对应的表达式元素加上量词，最终得到`[-+]?(\d+)?(\.\d+)?`。初看起来，这样并没有错，仔细看看却发现，这个表达式中所有元素都不是必须出现的，换句话说，这个表达式匹配成功时，可以不匹配任何文本！

简单去掉量词并不能解决问题，这样又会错过许多本应该匹配的文本。真正要做的，是理清表达式各个部分的关系，尤其是“可能出现”和“必须出现”之间的关系。

9.1.2 可能出现

与普通字符串处理相比，“可能出现”可以算作正则表达式最明显的特征，也是最常用的逻辑。虽然它看起来很直观，但细说起来，它其实分为两种情况：从元素外来看，元素可能出现也可能不出现，或者出现次数不确定；从元素内来看，元素可能表现为一种形态，也可能表现为另一种形态。

第一种情况需要用量词。在匹配双引号字符串的正则表达式中，两个双引号字符是必须出现的，它们之间的文本可能出现，也可能不出现，如果出现，长度没有限制，所以用`*`来限制；在



匹配 tag 的正则表达式中，<和>之内的 tag 名，必须包括至少一个字符，所以用+来限制。

第二种情况需要使用字符组或多选结构。如果各种可能形态都是单个字符，则使用字符组，比如上一节匹配数字字符串正则表达式中的[-+]，它说明“此处可能出现+号，也可能出现-号”。如果某一种可能形态不只是一个字符，比如 this 或 that，则应该使用多选分支(this|that)。

回过头来看数字字符串的匹配，根据上一节的分析，符号部分、整数部分、小数部分（包括小数点）都是可能出现的，但这种“可能”其实是需要细分的。

符号部分的“可能出现”其实是“可能出现也可能不出现的”；整数和小数部分的可能出现情况要复杂一点，需要研究可能出现的几种形态：只出现整数部分；只出现小数部分；整数部分和小数部分都出现。

只出现整数部分	\d+
只出现小数部分	\.\d+
整数部分和小数部分都出现	\d+\.\d+

使用多选结构将这三种形态统一起来，得到(\d+|\.\d+|\d+\.\d+)。最后得到的正则表达式就是[-+]?(\d+|\.\d+|\d+\.\d+)。

有些人会觉得(\d+|\.\d+|\d+\.\d+)麻烦，所以会把多选结构的各个分支合并，得到(\d+)?\.\?\d+。这两个正则表达式能匹配的文本的确相同，但我并不推荐这样写正则表达式，因为它的逻辑不够清晰：对表达式(\d+)?\.\?\d+来说，最后的\d+是必须出现的。如果文本只包含整数部分，比如3，\d+匹配的是整数部分，但如果文本包含小数部分，比如3.14，则\d+匹配的是小数部分的数字。表达式(\d+|\.\d+|\d+\.\d+)虽然麻烦一点，却是一目了然。

这个例子在第3章出现过，当时使用的表达式[-+]?(\d+|\.\d+|\d+\.\d+)仍然不够完美，它可能错误匹配-.14。此时，[-+]?中真正匹配的是-，(\d+|\.\d+|\d+\.\d+)中真正匹配的是\.\d+。

要解决这个问题，可以将符号-和+分情况对待：如果是+，则之后的表达式有三种可能：只有整数部分；有整数和小数部分只有小数部分，所以用表达式(\d+|\.\d+|\d+\.\d+)匹配；如果是-，则之后的表达式只有两种可能，不可能出现“只有小数部分”的情况，所以用表达式(\d+|\d+\.\d+)匹配。

综合这些情况，最终得到的表达式就是(?:\d+|\.\d+|\d+\.\d+)|-?(\d+|\d+\.\d+))。

9.1.3 不能出现

“不能出现”是正则表达式中最难处理的。

最简单的“不能出现”可以直接使用排除型字符组，它表示“此处必须出现一个字符，但不能是某些字符”。比如匹配双引号字符串，首尾两个双引号之间的字符，都不能出现双引号字符，用`[^"]`表示，上一节说到，这部分长度可以为零，应当使用量词`*`，所以整个表达式就是`"[^"]*"`。

再比如匹配 `html tag`，在`<`和`>`之间“不能出现”`>`，用`[^>]`表示，上一节说到应当使用量词`+`，所以整个表达式就是`<[^>]+>`，它在第 3 章也出现过，这个表达式仍然可能错误匹配`</>`。为解决这个问题，还需要细分两种可能：如果`<`之后是`/`，则还必须出现至少一个不为`>`的字符，比如`</a>`，应当使用表达式`/[^>]+`；如果`<`之后出现的字符不是`/`，则在这个字符之后，`>`之前，还可能出现`>`之外的字符，并且可能不出现，如果出现，长度没有限制，比如在`<a>`中，在`[^>]`匹配之后，`>`之前，就没有任何字符了，所以应当使用表达式`[^/][^>]*`。

不过，这只是最简单的情况；更复杂的情况是，在某个字符串中不容许出现某个字符串，比如 E-mail 地址中的用户名（username）就是如此。

在第 3 章曾经讲过电子邮件地址的匹配（☞ 38），使用`[-\w.]{0,64}`匹配邮件地址中的用户名，但是如果仔细阅读 E-mail 的规范<sup>1</sup>就会发现，点号和横线都不能出现在用户名的开头，也不能连续出现，也就是说，`-user`、`.user` 和 `John..Doe` 都不是合法的电子邮件地址。

点号和横线不能出现在开头的情况比较好满足，可以用`\w`匹配“非点号非横线字符”，既然整个用户名不能超过 64 个字符，那么之后的字符串长度不超过 63 个字符。综合起来，得到`[\w][\w.]{0,63}`。可是，不容许出现连续两个点号的要求则很难满足，所以下面集中讨论“不超过 63 个字符”部分的匹配。

要求不能出现两个连续点号，许多人的直观反应是`[^.][^.]`，所以整个表达式改为`[\w.]{0,63}[^.][^.][\w.]{0,63}`。但是这样行不通，原因有两点：首先，在`[^.][^.]`前后的两个`[\w.]{0,63}`能匹配的文本，总长度其实在 0~126 之间，无论我们如何修改量词，也不能在两个量词之间建立联系，保证总长度在 0~63 之内；其次也是最重要的，`[^.][^.]`的真正意思是“找到连续的两个非点号字符”，正则表达式匹配时会尽力满足这一要求，即便是类似 `123..456` 这样明显包含两个连续点号的文本，`[^.][^.]`仍然可以在其中找到四处匹配：`12`、`23`、`45`、`56`，同时左右两侧的`[-\w.]{0,63}`仍然可以成功匹配。

还有人想到的是效仿排除型字符组，用`(^\.\.)`来表示“不能出现两个连续点号”，遗憾的是，这样也行不通，因为正则表达式只有排除型字符组，没有“排除型括号”。

那么，不妨反过来思考：不能出现两个连续点号，意思是这段文本中的所有字符的右侧都不能是两个连续点号，用环视表示就是`[-\w.](?!\.\.)`，这样的字符最多有 63 个。所以，就得

<sup>1</sup> 请参考 [http://en.wikipedia.org/wiki/Email\\_address](http://en.wikipedia.org/wiki/Email_address)，虽然根据规范，username 中还可以包含`!#$%*/?^`{|}~`等字符，但在实际应用中，许多邮件服务提供商只支持字母、数字、下划线、横线、点号，所以暂不考虑其他字符。

到了表达式 `([-\w.](?!\\.\\.)){0,63}`。请注意，因为 `(?!\\.\\.)` 不是单个字符也不是字符组，所以用量词限时，必须使用括号将整个子表达式分为一组；再在表达式的最开头加上之前提到的 `\w`，保证第一个字符的正确性。从例 9-1 可以看到，这个表达式完全可以保证字符串不会以点号开头，不会包含连续点号，也不会超出规定长度。

### 例 9-1 使用环视实现“不能出现”的逻辑

```
#注意验证时要在首尾加上\A和\Z
usernameRegex = r"\A\w{0,63}\Z"
#合法的用户名
re.search(usernameRegex, "abc123_") != None      # => True
re.search(usernameRegex, "abc1-2.3_") != None     # => True
#包含连续点号
re.search(usernameRegex, "abc1-2..3_") != None    # => False
#开头字符不合法
re.search(usernameRegex, ".abc1-2.3_") != None   # => False
re.search(usernameRegex, "-abc1-2.3") != None     # => False
#长度超过限制
re.search(usernameRegex, "0"*65) != None          # => False
```

也可以更进一步，把“第一个字符不能是点号或横线”也用环视表达，整个表达式就是 `(?![-.])([-\w.](?!\\.\\.)){1,64}`。实际上，这个表达式确实更直观、更容易理解，只是注意量词的下限应当改为 1，因为用户名不能是空字符串。

总结一下：正则表达式中的“不能匹配”，最简单的情况可以用排除型字符组直接表示，但它只能表示“某个字符不能出现”；如果要表示“某个字符串不能出现”，一般都要用到否定环视，其逻辑是：在文本中的每个位置，都用环视否定“不能出现”的字符串（除此之外，还有另一种逻辑，下面讲解验证操作时会看到）。不过，一些比较古老的工具（比如 Apache 1.3，以及 Linux/UNIX 下的某些工具）并不支持否定环视，所以使用时必须留意。

## 9.2 正则表达式的常见操作

正则表达式执行的操作，可以粗略分为三大类：**匹配**、**替换**、**切分**。其中，匹配是最基本的操作——使用正则表达式无非是“用正则表达式匹配文本”，这种说法没错，但太笼统，细究起来，广义的“匹配”又可以分为两类：**提取**和**验证**。所以此处将常见操作分为四类，逐一详细讲解。

### 9.2.1 提取

数据提取，就是“用正则表达式遍历整个字符串，找出能够匹配的文本”，它主要用来提取所需要的数据，常见的任务有：找出文本中的电子邮件地址，找出 HTML 代码中的图片、超链

接……提取数据时主要注意三点：**完整、精确、效率**。

首先来看完整和精确，它们关系很紧密，所以集中讲解（见图 9-1）。完整是要保证正则表达式能匹配足够多的文本，兼容可能出现的各种情况，避免疏漏；精确是则要保证正则表达式匹配足够少的文本，避免错误匹配其他文本。

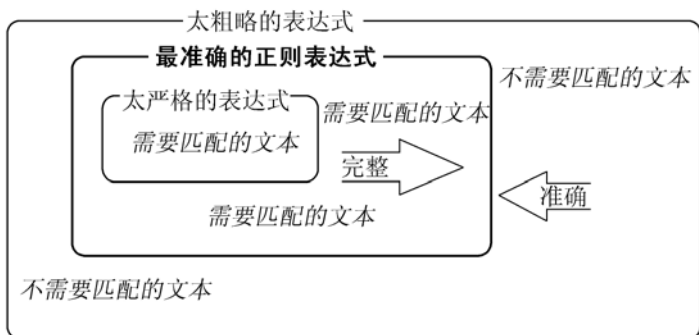


图 9-1 完整和精确

仍然举之前匹配浮点数的例子，回顾逐步完善这个表达式的过程。

(1) 常见的浮点数字符串就是 **3.24**、**0.618** 之类，也就是“数字字符串 + 小数点 + 数字字符串”，所以用表达式 `\d+\.\d+`。

(2) 它无法匹配只有小数点和小数部分的浮点数，比如 **.618**，所以应当把整数部分改为可能出现的，得到 `\d*\.\d+`。

(3) 它无法匹配只有整数部分的浮点数，比如 **3**，所以应当把小数点和小数部分改为可能出现的，得到 `\d*\.\?\d*`。

(4) 表达式中的小数点和小数部分是一体的，要么同时出现，要么同时不出现，所以应当用括号分组功能，得到 `\d*(\.\d)*`。

(5) 表达式中，所有元素都有量词限制，可以不匹配任何字符，所以它可以匹配空字符串，或者匹配单独的点号，所以应当分情况考虑，得到 `(\d+|\.\d+|\d+\.\d+)`。

(6) 表达式还不能匹配包含符号的浮点数，比如 **+0.8**、**-0.7** 之类，所以还应当加上对符号的处理，得到 `[+]?(\d+|\.\d+|\d+\.\d+)`。

(7) 这个表达式可能匹配 **-7**，但通常它并不是一个合法的浮点数字符串，而应当写作 **-0.7**，所以应当再次修改，最后得到的表达式就是 `(+?(\d+|\.\d+|\d+\.\d+)|-?(\d+|\d+\.\d+))`。

其中第 2、第 3、第 6 步，都是为了兼容更多情况而做的修改，也就是保证匹配的完整性。但是，这种修改可能有副作用，影响整个表达式的匹配，所以第 4、第 5、第 7 步是给这些修改打上补丁，保证准确性，避免错误匹配其他文本。

不过，这个表达式还不够精确——因为提取操作往往是在大段文本中进行的，真正感兴趣的文本并不是独立存在的——如果遇到日期字符串 `2010.12.22`，就会错误提取出 `2010.12`；同样，如果用表达式 `\d{0,6}` 提取邮政编码，遇到手机号码 `13801234567`，也可能错误提取出 `138012`。

为解决这个问题，保证精确性，需要**为正则表达式设定合适的断言**。针对浮点数的匹配，为了避免日期字符串的错误匹配，可以在表达式首尾加上否定环视 `(?!<\.)` 和 `(?!\.)`，保证匹配文本的两端不能为点号，于是得到正则表达式 `(?!<\.)(+?(\\d+|\\.\\d+|\\d+\\.\\d+)|-?(\\d+|\\d+\\.\\d+))(?!\\.)`。这样，就真正避免了 `2010.12.22` 的错误匹配。

提取数据是正则表达式常见的操作，而完整和精确又是非常重要的方面，所以值得多花点篇幅，下面再举一个匹配双引号字符串的例子。

前面用来匹配双引号字符串的正则表达式 `"[^"]*"` 其实并不够完整——因为双引号字符串内部还可能嵌套转义的双引号。如果真的出现这样的双引号，就可能出现意外：

字符串

```
"quoted string with \"nested quote\""
```

匹配结果

```
"quoted string with \"
\""
```

因为字符串中出现了两个转义序列 `\"`，以 `"[^"]*"` 来匹配就会分开匹配两次。

改进这个正则表达式，最容易想到的是使用逆序环视，约定 `"` 之前必须出现反斜线字符 `\`，所以 `"[^"]*"` 要改为 `"((?<=\\)"|"[^"]")"`，因为 `(?<=\\)"` 能匹配字符串中的 `\"`，所以 `"quoted string with \"nested quote\""` 就可以完整匹配。

现在，完整的问题已经解决。下面来看精确，尝试用这个表达式匹配更复杂的文本，看看是否会错误匹配不应该匹配的文本。

字符串

```
"escaped backslash\\" and "another"
```

匹配结果

```
"escaped backslash\\" and "
```

字符串中只出现了一个转义序列 `\`，表示反斜线字符。因为 `((?<=\\)"|"[^"]")` 中的 `(?<=\\)"` 匹配字符串中的 `\`，没有准确识别字符串中的转义序列 `\`。

看来应该修改逆序环视: `((?<=\\)"|^[^"])` 匹配的 `\` 之前不能再出现 `\`, 也就是说 `"` 之前应当出现 `\`, 但不能出现 `\\`, 所以需要用嵌套环视把表达式改为 `"((?<!\\\\)(?<=\\)"|"[^"]")"`; 这样, 就解决了 `"escaped backslash\\"` and `"` 匹配的问题。

不过还有更复杂的情况, 比如 `"double escaped backslash\\\\"` 是完整的双引号字符串, 其中出现了相连的转义序列: `\\`, 而 `(?<!\\\\)(?<=\\)"` 要求 `"` 之前不能出现两个连续的 `\`, 所以会出错。解决思路之一是, 约定 `"` 之前出现的 `\` 的数目为奇数, 可是除去 .NET, 大多数语言中的逆序环视都要求表达式匹配的文本长度是固定的, 所以这个办法并不实用。那么, 有没有通用的解法呢?

如果从处理一般的转义序列出发, 用 `\\.` 匹配转义序列, 就可以一次匹配整个转义序列, 而不只是匹配转义中 `\` 之后的字符, 把表达式改为 `"(\\.|[^\"])*"`。这时候, 既可以完整匹配 `"escaped backslash\\\\"` and `"another"`, 也可以匹配 `"escaped backslash\\"` and `"another"` 中的 `"escaped backslash\\"` 和 `"another"`。

现在完整性的问题基本解决, 为确保精确性还必须确认错误的字符串不会匹配。比如字符串 `"quoted string without ending \"`, 除去开头的双引号和转义的双引号, 没有结尾的双引号, 所以不是一个合法的双引号字符串, 也就不应该能匹配, 但可以被 `"(\\.|[^\"])*"` 错误匹配。

这是为什么呢? 来分析匹配的过程: 字符串结尾的 `\`, 本来应该由多选结构中的 `\\.` 匹配, 这时候字符串中已经没有剩下字符了, 而表达式结尾的 `"` 还没有匹配, 所以需要回溯。这时候, `\\.` “交还” 的 `\` 中, `\` 可以由 `[^"]` 匹配, 还剩下 `"`, 正则表达式中也剩下 `"` 没有匹配, 所以匹配完成——尽管结果是错的。

问题的症结在于, 多选结构 `(\\.|[^\"])` 中, `\\.` 交还的部分可能被 `[^"]` 匹配, 而我们本意是希望这两个多选分支应对不同的情况, 前者匹配转义序列, 后者匹配非转义序列。但是仔细想想就能知道, 在双引号字符串中的非转义序列中, 不但不能出现 `"`, 也不能出现 `\`。讲解多选结构时曾经提到: 应当尽可能避免多选分支匹配相同的文本, 所以应当把多选结构改为 `(\\.|[^\\""])`, 整个表达式改为 `"(\\.|[^\\""])*"`。

最后得到的正则表达式才算是真正保证了完整性——可以匹配包含转义双引号的双引号字符串, 也保证了精确性——包含转义双引号, 但没有结束双引号的非法双引号字符串, 是不能匹配的。

关于完整和准确, 最后看一个例子。第 3 章提到用正则表达式提取超链接, 当时使用的正则表达式是 `<a\s+href\s*=\s*["']?(["'\s]*+)[["']]?>[<+]/></a>`, 其中 `(["'\s]*+)` 用来匹配超链接的地址, `[<+]` 用来匹配超链接的文本。但是, 这个表达式其实只能匹配 `<a href="url">text</a>` 这类“标准”的超链接。

然而在 `<a` 和 `href` 之间, 可能还存在有其他内容, 比如 `<a id="mh_img" href...`。为了应付这样的情况, 似乎可以将 `<a\s+href` 中换成 `<a\s[>]*\s?href`, 这样既能匹配 `<a href...`, 又能

匹配<a id="mh_img"></a>[^>]\*\s?中的`[^>]`和`\s`都不是必须出现的，所以整个表达式也可以匹配

仔细想想，`<a`和`href`之间出现其他内容，可以由表达式`[^>]*`匹配；同时要保证`<a`之后、`href`之前都出现空白字符，最合适使用可以进行“原地判断”的环视，`<a(?:=\s)`和`(?<=\s)href`匹配了符合条件的`<a`和`href`，同时省去了必须出现两个的`\s`的限制，所以匹配这一段的表达式就是`<a(?:=\s)[^>]*(?<=\s)href`。

再看`href`属性之后的部分，之前的表达式中用`[\\"']?>`来匹配，但也不够完整，因为其中还可能其他内容，比如`<a href="someurl" data="somedata">`，不过这种情况处理起来很简单，用`[^>]*`就可以处理，所以`[\\"']?>`要改为`[\\"']?[^>]*>`。

最后考虑`<a>`和`</a>`之间的内容，之前用`([^<]+)`匹配（添加括号分组是为了提取出这段内容），是考虑这里的内容不包含其他tag，比如`<a href="someurl">text</a>`，但实际情况并非如此，其中可能是图片`<a href="someurl"></a>`。因为这段内容可能包括任何字符，甚至包括换行符，所以最好用`[\s\S]`（如果指定了单行模式，也可以使用点号）匹配其中的字符，比较好的办法是使用忽略优先量词，以`([\s\S]+?)`匹配，忽略优先量词保证了这个表达式的匹配到`</a>`之前结束。

综合起来，兼顾完整和准确的匹配超链接的表达式就是`<a(?:=\s)[^>]*(?<=\s)href\s*=\s*[\\"']?(?:[\\"'\s]+)[\\"']?[^>]*>([\s\S]+?)</a>`。

完整和准确很重要，也很费心思；但也有很多时候不用考虑这么细致，这取决于你要处理的文本。如果要处理的文本中只会出现像 3.14 一样规则的浮点数，完全可以放心使用`\d+.\d+`；如果确认文本中不会出现手机号码之类的字符串，完全可以放心使用`\d{0,6}`；如果要处理的都是简单的双引号字符串，完全可以放心使用`"[^"]*".....`。

提取数据时，应当注意的另一点是效率。在讲解正则表达式的匹配原理时已经提到过，在编写表达式时需要注意效率；这里要说的是另一方面：提取的效率不仅与正则表达式本身有关，也与调用的 API 有关。许多网站每天的访问日志经常有到几百兆甚至更大，如果先读入整个日志文件，再一次性提取某些数据（比如访问者 IP），就可能消耗非常长的时间。

遇到这种情况，可取的办法是在条件容许的情况下（比如只需要逐步提取出来，依次处理），逐步进行提取——提取一个结果，处理，再提取一个结果，再处理……表 9-2 列出了常用语言中的提取 API。

表 9-2 常用语言中的提取 API 方法

语言	方法	备注
.NET	<code>Regex.Match(string)</code>	逐次进行
	<code>Regex.Matches(string, regex)</code>	一次性进行

(续表)

语言	方法	备注
Java	<code>Matcher.find()</code>	逐步进行
JavaScript	<code>RegExp.exec(string)</code>	一次性进行
	<code>string.match(RegExp)</code>	一次性进行
PHP	<code>preg_match(pattern, subject, matches)</code>	逐步进行
	<code>preg_match_all(pattern, subject, matches)</code>	一次性进行
Python	<code>re.find(pattern, string)</code>	逐步进行
	<code>re.finditer(pattern, string)</code>	逐步进行
Ruby	<code>Regexp.match(str)</code>	只能找到第一次匹配
	<code>string.index(Regexp, int)</code>	逐步进行
	<code>string.scan(Regexp)</code>	一次性进行

一次性提取所有匹配结果的操作不必说太多；这里只讲解“逐步进行”时如何真正保证“逐步”？也就是说，在第二次调用匹配时，如何保证是“承接”第一次调用，找到下一个匹配结果。下面介绍三种常用的办法，例子是查找字符串 123 45 6 中的数字字符串，依次输出 123、45、6。

如果使用函数式正则表达式采用的是面向对象式处理，匹配结果对象可以自动保存匹配状态信息，其中包括匹配本次匹配结束的位置，所以下次匹配时，会自动从之前记录的结束位置开始。Java 就是这样，从例 9-2 可以看到，循环调用 `Matcher.find()` 方法，就可以逐个获得所有匹配（在.NET 中，是循环调用 `Match.NextMatch()`）。

例 9-2 面向对象式处理中的逐次提取(Java)

```
String str = "123 45 6";
Pattern p = Pattern.compile("\\d+");
Matcher m = p.matcher(str);
while (m.find()) {
    System.out.println(m.group());
}
```

如果使用函数式处理，则必须自己保存匹配的状态信息，手动指定偏移值。多数语言都有办法在匹配时指定偏移值，也就是“从字符串的 *offset* 位置开始尝试匹配”。如果要逐一获得所有匹配，每次匹配时手动记录上次匹配的结束位置，下次匹配时，指定从此位置开始尝试即可<sup>1</sup>。

<sup>1</sup> 注意，字符串处理时可能有人习惯将偏移值指定为“上一次匹配的起始位置+1”，但正则表达式处理时这样是不对的，比如正则表达式是 `\d+`，而字符串是“123 45 6”，第一次匹配的结果是 123，如果把偏移值设定为“上一次匹配的起始位置+1”，之后的匹配结果就是 23，3……



在 PHP、JavaScript、Ruby、Python 中，通常采用这种办法，例 9-3 给出了 PHP 的代码，请注意其中的 *lastOffset* 变量。

#### 例 9-3 面向对象式处理中的逐次提取(PHP)

```
$string="123 45 6";
$regex="/\\d+/";
$matched=1;
$oneMatch=array();
$lastOffset=0;
$matched = preg_match($regex, $string, $oneMatch, PREG_OFFSET_CAPTURE, $lastOffset);
while ($matched == 1) {
    //上次匹配开始位置+上次匹配文本的长度，得到上次匹配的结束位置
    $lastOffset = $oneMatch[0][1] + strlen($oneMatch[0][0]);
    echo $oneMatch[0][0]."<br />";
    $matched = preg_match($regex, $string, $oneMatch, PREG_OFFSET_CAPTURE, $lastOffset);
}
```

最后一种办法是使用迭代器。例 9-4 展示了在 Python 中，调用 `re.finditer(pattern, string)`，会得到一个迭代器（虽然看起来是用 for 循环，但每次的匹配结果其实是临时计算得到的），这种办法目前只有 Python 提供。

#### 例 9-4 使用迭代器进行逐次提取(Python)

```
for match in re.finditer(r"\d+", "123 45 6"):
    print match.group(0),
```

## 9.2.2 验证

数据验证是另一类“匹配”操作，它用来“检查字符串能否完全由正则表达式匹配”，主要用来测试和保证数据的合法性，比如 Web 编程中经常用正则表达式验证用户输入数据，比如手机号码、邮件地址、QQ 号码等等。

初看起来，这也是用正则表达式在字符串中查找匹配文本。仔细分析又不一样：提取返回的是字符串，验证返回的是布尔值；提取时，正则表达式最终匹配的起始/结束位置都是不确定的，需要逐次试错才能确定；验证时，正则表达式的起始/结束位置必须定位于字符串的起始/结束位置，从字符串起始位置开始尝试匹配，只要表达式中任意一个必须匹配的元素出现无法匹配的情况，即宣告验证失败，不必逐次试错。假设正则表达式是 `\d{6}`，字符串是 `a123456b`，很明显，提取操作中应当返回 `123456`，但验证操作应当返回 `False`，因为表达式匹配的并不是整个字符串。

针对验证操作的特点，有些语言中提供了专门的方法进行正则表达式验证，它们保证正则表达式匹配的起始/结束位置就是字符串的起始/结束位置，而且返回布尔值。如果没有提供专门的方法，通常的办法是手动在正则表达式的首尾加上匹配字符串起始/结束位置的 `\A` 和 `\Z`，再通过

返回值判断<sup>1</sup>。

这里，有必要单独谈谈`\A`和`\Z`的问题。许多示范代码（包括本书开头部分）都使用`^`和`$`匹配字符串的起始/结束位置，但第 63 页讲过，`$`匹配的其实是“结尾换行符之前的位置”，如果字符串的结尾是换行符`\n`，即便在表达式末尾添加了`$`，匹配时仍然会忽略最后的换行符。所以应当使用真正匹配“结束位置”的锚点`\Z`（Java、JavaScript、.NET、Ruby、PHP 都支持`\Z`，只有 Python 不支持，Python 中的`\Z`等于其他语言中的`\z`）。虽然`^`和`\A`都是匹配字符串的起始位置，但为了与`\Z`对应，推荐使用`\A`。而且使用`\A`和`\Z`还有一个好处：如果设定了多行模式，`^`和`$`的匹配会有所不同，`\A`和`\Z`则不受影响。表 9-3 中列出了常用语言中的验证方法。

表 9-3 常用语言中的验证方法

语言	验证方法	备注
.NET	<code>Regex.IsMatch(string, regex)</code>	返回布尔值，但需要 <code>\A</code> 和 <code>\Z</code>
Java	<code>String.matches(regex)</code>	专用于验证，返回布尔值，不需要 <code>\A</code> 和 <code>\Z</code>
JavaScript	<code>RegExp.test(string)</code>	返回布尔值，需要 <code>\A</code> 和 <code>\Z</code>
PHP	<code>preg_match(pattern, subject)!=0</code>	需要 <code>\A</code> 和 <code>\Z</code>
Python	<code>re.search(pattern, string)!=None</code>	需要 <code>\A</code> 和 <code>\Z</code>
	<code>re.match(pattern, string)!=None</code>	不需要 <code>\A</code> ，只需要 <code>\Z</code>
Ruby	<code>Regexp.match(string)!=nil</code>	需要 <code>\A</code> 和 <code>\Z</code>
	<code>String.index(Regexp)!=nil</code>	需要 <code>\A</code> 和 <code>\Z</code>

除去方法上的差别，验证与匹配在正则表达式编写上也有区别。因为验证时文本的起始/结束位置是预先知道的，所以验证的表达式编写起来更加简单。比如之前匹配浮点数的表达式，首先得到的正则表达式是`(+?(\\d+|\\.\\d+|\\d+\\.\\d+)|-?(\\d+|\\d+\\.\\d+))`，进行数据提取还需要在两端加上环视，防止错误匹配其他字符；但是如果是数据验证，就可以完全忽略考虑两端的环视，应该/不应该出现什么字符，直接在首尾加上`\A`和`\Z`即可，最终，验证浮点数的表达式是`\A(+?(\\d+|\\.\\d+|\\d+\\.\\d+)|-?(\\d+|\\d+\\.\\d+))\Z`。同样道理，匹配邮政编码的正则表达式是`\\d{6}`，如果用于数据验证，可以直接使用`\A\\d{6}\Z`。

其实，根据这两点差别，完全可以为验证操作找到完全不同的思路：从根本上说，数据提取的结果是字符串，而验证的结果是布尔值。既然返回布尔值，就很容易想到布尔运算，所以验证时不妨简单将各个条件叠加起来，直接得到最后的表达式。下面以密码字符串的验证为例，说明这种思路。

用正则表达式验证密码字符串是否合法，“合法”密码必须同时满足以下所有要求。

<sup>1</sup> 专用于验证的方法应当只返回布尔值，所以表 9-3 中对原本不返回布尔值的函数进行了包装，以返回布尔值。

- (1) 密码长度在 6~12 个字符之间；
- (2) 只能由小写字母、阿拉伯数字、横线组成；
- (3) 开头和结尾不能是横线；
- (4) 不能全部是数字；
- (5) 不容许有连续（2 个及以上）的横线。

针对这 5 点要求，可以分别列出对应的表达式。

- (1) 密码长度在 6~12 个字符之间：其形式类似 `{6,12}`；
- (2) 只能由小写字母、阿拉伯数字、横线组成：所有的字符都只能由 `[0-9a-z-]` 匹配；
- (3) 开头和结尾不能是横线：开头 `\A(?!-)`，结尾 `(?<!--)\z`；
- (4) 不能全部是数字：也就是说，必须出现一个 `[^0-9]` 或者 `\D`；
- (5) 不容许有连续（2 个及以上）的横线：也就是说，不能出现 `--`。

如果用来提取数据，就必须使用能兼顾这 5 点要求的表达式，也就是把上面 5 个表达式综合起来。前三点比较好处理，可以合并为 `\A(?!-)[0-9a-z-]{6,12}(?<!--)\z`，剩下的两点要求就麻烦多了。如果你仔细阅读过前面关于三种逻辑的讲解，或许知道能用环视结构解决，而且在验证中完全可以更简单，只依靠环视解决所有的问题。

既然验证只需要返回布尔值，不妨换个思路，针对每个要求分别判断，再将各个要求的布尔值取“与”（and）运算即可。取得布尔值，最容易想到的就是环视功能，所以下面列出针对 5 点要求的环视。其中值得注意的是第 5 点“不能出现”的处理，之前我们用“每个位置都不能出现”的逻辑来满足“不能出现”的要求，但是此处用到了另一种逻辑：先用子表达式 `\A.*--` 匹配“不应该出现”的文本，再用否定环视来杜绝这种情况。

- (1) 密码长度在 6~12 个字符之间： `(?=\A.{6,12}\z)`；
- (2) 只能由小写字母、阿拉伯数字、横线组成： `(?=\A[0-9a-z-]*\z)`；
- (3) 开头和结尾不能是横线：开头 `(?!\\A-)`，结尾 `(?!.*\\z)`；
- (4) 不能全部是数字： `(?=\A.*[^0-9])`（这里不需要出现 `\z`，只要出现了非数字字符就可以）；
- (5) 不容许有连续（2 个及以上）的横线： `(?!\\A.*--)`（这里不需要出现 `\z`）。

第 4 章提到过，环视是在原地进行的，而且可以组合起来。而根据环视的组合规则，简单并列多个环视的意思是：在当前位置，所有环视的判断都必须成功。不妨直接选择字符串的起始位置作为“当前位置”，用 `\A` 匹配，再将所有环视结构列在后面即可。例 9-5 说明，这个表达式确实可以满足要求。

## 例 9-5 纯粹使用环视完成验证

```
#Python 不支持\z, 所以\z 必须写成\Z
passwdRegex = r"\A(?:\A.{6,12}\Z)(?=\A[0-9a-z-]*\Z)(?! \A-)(?!.*-\Z)(?=\A.*[^\0-9])
(?:\A.*--)"
#合法
re.search(passwdRegex, "abc123-d") != None           # => True
#超长
re.search(passwdRegex, "abcdefg1234567890") != None   # => False
#出现了其他字符
re.search(passwdRegex, "1234ab*") != None             # => False
#开头或结尾出现了横线
re.search(passwdRegex, "abc123-") != None             # => False
re.search(passwdRegex, "-abc123") != None             # => False
#出现了连续横线
re.search(passwdRegex, "abc--123") != None            # => False
```

如果你觉得这个表达式有点复杂, 可以对它做一些整理: 开头的`\A`定位在字符串的起始位置, 而环视不会真正改变匹配的位置, 每个环视判断完成, 下一个环视仍然从字符串的起始位置开始, 因此可以将各个环视中的`\A`去掉, 但是末尾的`\z`不能去掉, 否则`(?=.{6,12})`就不能保证整个字符串的长度在 6~12 之间。这个表达式的测试结果见例 9-6。

## 例 9-6 用修改后的表达式完成验证

```
#Python 不支持\z, 所以\z 必须写成\Z
passwdRegex = r"\A(?:=.{6,12}\Z)(?=[0-9a-z-]*\Z)(?!-)(?!.*-\Z)(?=\A.*[^\0-9])
(?:\A.*--)"
#合法
re.search(passwdRegex, "abc123-d") != None           # => True
#超长
re.search(passwdRegex, "abcdefg1234567890") != None   # => False
#出现了其他字符
re.search(passwdRegex, "1234ab*") != None             # => False
#开头或结尾出现了横线
re.search(passwdRegex, "abc123-") != None             # => False
re.search(passwdRegex, "-abc123") != None             # => False
#出现了连续横线
re.search(passwdRegex, "abc--123") != None            # => False
```

这种思路简单易行, 对于数据验证特别有用, 因为验证操作一般会明确给出验证规则(比如这个例子中的五点要求)。相比之下, 提取操作不一定会给出关于要提取文本的严格描述, 所以一般处理是先找到最明显特征, 再逐步完善表达式。

9.2.3 替换

替换操作也是一种“匹配”，它通常需要三个参数，其中 *regex*（或 *pattern*）表示查找需要替换文本的正则表达式，*replacement* 表示“要替换成”的字符串，*subject*（或 *string*、*input*）则是要进行替换处理的文本。表 9-4 中列出了常用语言中的替换方法。

表 9-4 常用语言中的替换方法

语言	方法	备注
.NET	<code>Regex.Replace(input, pattern, replacement)</code>	静态方法
	<code>Regex.Replace(input, replacement)</code>	<code>Regex</code> 实例事先生成
Java	<code>String.replaceAll(regex, replacement)</code>	静态方法
	<code>String.replaceFirst(regex, replacement)</code>	同上，但只进行第一次替换
	<code>Matcher.replaceAll(replacement)</code>	<code>Matcher</code> 实例事先根据 <code>Pattern</code> 和 <code>String</code> 生成
	<code>Matcher.replaceFirst(replacement)</code>	同上，但只进行第一次替换
JavaScript	<code>string.replace(RegExp, replacement)</code>	<i>replacement</i> 为普通函数
	<code>string.replace(RegExp, replacement)</code>	<i>replacement</i> 为回调函数
PHP	<code>preg_replace(pattern, replacement, subject)</code>	
	<code>preg_replace_callback(pattern, callback, subject)</code>	<i>callback</i> 为回调函数
Python	<code>re.sub(pattern, replacement, string)</code>	
	<code>re.sub(pattern, replacement, string)</code>	<i>replacement</i> 为回调函数
Ruby	<code>String.sub(Regexp, String)</code>	只替换第一次匹配
	<code>String.gsub(Regexp, String)</code>	替换所有匹配

注：以上列出的都是最基本的形式，在默认情况下，替换是对所有的匹配实行的，也就是说，匹配发生了 *n* 次，替换就会发生 *n* 次，但替换发生的最大次数往往可以用一个可选参数指定，如果这个参数设定为负数，则会对所有匹配进行替换。具体细节此处不展开，请参考各语言的详细章节。

最简单的替换是删除不需要的文本，也就是将 *replacement* 设定为空字符串“”，此时 *regex* 能匹配的所有文本替换为空字符串，也就是“删除”了。比如将正则表达式 `\s+` 能匹配的所有文本都替换为空字符串，就删除了所有的空白字符（如果要处理的文本中包含 Unicode 空白字符，比如中文的全角空格，推荐使用 Unicode Property `\p{Z}`，(🔗 123)。当然有时候情况没有这么简单，第 71 页的正则表达式可以处理中英文混排的文本，只删除中文文本中的空白字符，而保留英文文本中的空白字符。

替换也可以用来调整数据的格式。比如日期字符串的整理，如 2010 年 12 月 22 日，按照

MM/DD/YYYY 格式就是 12/22/2010，按照 DD/MM/YYYY 格式就是 22/12/2010，按照 YYYY-MM-DD 格式就是 2010-12-22。使用正则表达式替换很容易进行这种调整，具体的例子在第 44 页已经讲解过，这里不再重复。

关于这类替换，唯一要注意的是，如果在 *replacement* 中出现了对捕获分组的引用，同样需要转义，因为 *replacement* 是字符串变量，而 `\1`、`\2` 之类是正则文字中的记法，不是字符串中的记法，代码如例 9-7 所示。

例 9-7 replacement 中的引用也需要转义

```
print re.sub(r"(\d{4})-(\d{2})-(\d{2})", r"\2/\3/\1", "2010-12-22")
12/22/2010

print re.sub(r"(\d{4})-(\d{2})-(\d{2})", "\\2/\\3/\\1", "2010-12-22")
12/22/2010

#因为\2、\3、\1 表示 ASCII 码值为 2、3、1 的字符，无法显示，所以用 repr()展示
print repr(re.sub(r"(\d{4})-(\d{2})-(\d{2})", "\2/\3/\1", "2010-12-22"))
'\x02/\x03/\x01'
```

表 9-5 和表 9-6 分别总结了常用语言中对普通分组和命名分组的引用记法。

表 9-5 常用语言中对分组的引用

语言	表达式中的反向引用	替换中的反向引用
.NET	<code>\num</code>	<code>\$num</code>
Java	<code>\num</code>	<code>\$num</code>
JavaScript	<code>\num</code>	<code>\$num</code>
PHP	<code>\num</code>	<code>\num</code> 或 <code>\$num</code> (PHP 4.0.4 以上版本)
Python	<code>\num</code>	<code>\num</code>
Ruby	<code>\num</code>	<code>\num</code>

表 9-6 常用语言中对命名分组的引用

语言	分组记法	表达式中的引用记法	替换时的引用记法
.NET	<code>(?&lt;name&gt;...)</code>	<code>\k&lt;name&gt;</code>	<code>\${name}</code>
PHP	<code>(?P&lt;name&gt;...)</code>	<code>(?P=name)</code>	不支持，只能使用 <code>\\${num}</code> ，其中 <code>num</code> 为对应分组的数字编号
Python	<code>(?P&lt;name&gt;...)</code>	<code>(?P=name)</code>	<code>\g&lt;name&gt;</code>
Ruby	<code>(?&lt;name&gt;...)</code>	<code>\k&lt;name&gt;</code>	<code>\k&lt;name&gt;</code>

还有些时候，虽然想进行的操作是替换，但是不能只关心需要被替换的那段文本，而必须找到

更多的文本，替换其中一部分，同时需要留下的部分。这样看起来有点多此一举，其实却是必需的，比如删除代码文件行末的`//comment` 注释就是如此：在许多编程语言中，都可以在行末添加`//comment` 注释，它以`//`开头，中间可以出现任何字符，直到该行结束为止。

替换前

```
String s = "some text"; //initialization
```

替换后

```
String s = "some text";
```

要去掉末尾的注释，最容易想到的就是用`//.*$`匹配注释文本（必须设定多行模式，否则`$`无法匹配文本内部的行结束位置），将它替换为空字符串。这种做法处理上面的情况当然没问题，但遇到有些情况则会出错。

替换前

```
String url = "http://somehost.com"; //initialization
```

替换后

```
String url = "http:
```

解决办法似乎是要给正则表达式`//.*$`加上限制，在`//`匹配之后、`$`匹配之前，不能再出现`/`，所以把表达式改为`//[^/]*$`。这样解决了上面的问题，但是在有些情况下还是会出错。

替换前

```
String url = "http://somehost.com";
```

替换后

```
String url = "http:
```

或许可以再加点限制，在`//`匹配之后、`$`匹配之前，不但不能出现`/`，也不能出现`"`，所以把表达式改为`//[^/"]*$`，就解决了上面的问题。不幸的是，这个表达式还有问题，有些情况下，替换前后不会发生任何变化。

替换前

```
String s = "some text"; //initialization with "some text"
```

替换后

```
String s = "some text"; //initialization with "some text"
```

仔细思考就会发现，虽然注释部分是从`//`开头，到`$`结尾，但为了准确判断它，不能只着眼于注释本身，而必须考虑更多：粗略来说，`//`确实是注释的开头，但它也有可能出现在双引号字符串内部，或者说，在真正的程序代码中。遇到这样的`//`，应当忽略——或者说，应当忽略双引号字符串内部的`//`。

那么换个角度，从整行文本的开头来看：首先是可能出现的真正的程序代码，如果`//`出现在之中，肯定是出现在双引号字符串内部（暂时不考虑某些语言中单引号字符串的情况）；之后，如果出现了`//`，则肯定是注释部分的开始。

可以将真正程序代码的文本分为三种情况：单独出现的`/`、双引号字符串、`/`和`"`之外的文本，

将对应的三个表达式综合起来，就得到匹配真正程序代码的表达式，再加上匹配末尾的 `///.*$`，就得到了分段匹配整行的正则表达式，整个表达式的匹配过程如表 9-7 所示。

表 9-7 匹配包含注释的代码行的正则表达式

部分		表达式
程序	单独出现的/	<code>(?&lt;!/)/(?!/)</code>
	双引号字符串	<code>"(\\. [^\\""])*"</code>
	/和"之外的文本	<code>[^/"]</code>
注释		<code>///<code>.*\$</code></code>
综合		<code>^((?&lt;!/)/(?!/) "(\\. [^\\""])*" [^/"])*///<code>.*\$</code></code>

注 1：虽然源代码中并不是每一行之后都会有注释，但在这个表达式中，匹配注释部分的 `///.*`并没有用量词限定——如果这一行不包含注释，则根本不需要进行替换处理，所以正则表达式无法匹配，反而是提高了效率。

注 2：必须指定多行模式，同时用 `^`和`$`来限定这个表达式只匹配一行文本。

这个表达式匹配了整行的文本，为去掉末尾的注释，必须决定保留下哪些内容。在表达式中，匹配程序语句的部分是 `^((?<!/)/(?!/)|"(\\.|[^\\""])*"|[^/"])*`，在第 47 页提到，如果某个捕获分组使用了 `*`量词，则对应分组保留的是最后一次匹配时捕获的文本，所以如果引用这个分组匹配的文本，得到的只是其中某个子表达式匹配的文本，解决办法是给这个表达式加上括号，这样编号为 1 的分组就保存了整个程序语句了。从例 9-8 可以看到，这个表达式确实可以准确删除所有的注释。

例 9-8 准确删除注释

```
wholelineRegex = r"^(?m)^(((?<!/)/(?!/)|\"(\\.|[^\\""])*\"|\"[^\"]\")*)///.*$\"
print re.sub(wholelineRegex, r\"\\1\", \"String s = \\\"some text\\\"; //initialization\\\")
String s = \"some text\";

print re.sub(wholelineRegex, r\"\\1\", \"String url = \\\"http://somehost.com\\\";
//initialization\")
String url = \"http://somehost.com\";

print re.sub(wholelineRegex, r\"\\1\", \"String url = \\\"http://somehost.com\\\";\")
String url = \"http://somehost.com\";

print re.sub(wholelineRegex, r\"\\1\", \"String s = \\\"some text\\\"; //initialization with
\\\"some text\\\")
String s = \"some text\";
```

一般来说，`replacement` 字符串足够表示替换的结果了。不过也有时候需要进行更复杂的替换，所以在有些语言的正则表达式替换方法中，`replacement` 可以是字符串，也可以是一个回调



函数（一般记为 *callback*），比如 Python 就是如此，下面举两个例子。

假设文本中包含一些浮点数，其格式都是 2.618、10.23 之类，不包含符号，而且整数部分、小数点必定出现，现在希望处理所有的浮点数，只保留一位小数，而且做四舍五入处理，也就是说，2.618 要变成 3.0，而 10.23 要变为 10.2。

匹配浮点数的正则表达式很简单，粗略记为 `\d+\.\d+`，但无论如何编写 *replacement* 字符串都无法进行四舍五入处理；如果使用 *replacement* 函数，则非常简单，它接收的是一个 MatchObject 对象，返回一个字符串。程序代码如例 9-9 所示。

#### 例 9-9 使用 replacement 函数进行四舍五入处理

```
#先定义四舍五入处理的函数
def processFloat(matchObj) :
    matchedStr = matchObj.group(0)
    #小数点之后数字为 01234 的情况
    if re.search("\.[0-4]", matchedStr) != None :
        #保留整数部分和小数第一位
        return "%.1f" % float(matchedStr)
    else :
        #其他情况，取整数部分+1，小数第一位补 0
        return "%.1f" % (int(float(matchedStr))+1)

print re.sub(r"\d+\.\d+", processFloat, "2.618 10.236")
3.0 10.2
```

还有更复杂的情况：有些网页为了防止抓取，在源代码中将文字都用 `&#dec;` 或者 `&#xhex;`；这种形式的 HTML Entity 表示，其中的 *dec* 或 *hex* 分别是十进制和十六进制编码值。用户在浏览器中看到的“收发”两个字，对应的 HTML 代码可能是 `&#x6536;&#x53d1;` 或者 `&#25910;&#21457;`，其中 6536 和 25901 分别是“收”的 Unicode 码值的十进制和十六进制编码值，而 53d1 和 21457 分别是“发”的 Unicode 码值的十进制和十六进制编码值。为了获取网页的内容，需要将 HTML Entity 还原为对应的文字。对于 `&#dec;` 或者 `&#xhex;` 这类文本的匹配，正则表达式很擅长，但是将某个编码值转换为对应的文字正则表达式却是无能为力的。这种时候，用 *replacement* 函数处理起来却非常简单：正则表达式是 `&#(\d+|x[0-9a-fA-F]+);`，`&#` 是必须出现的，然后是编码值字符串，如果是十进制的，则用 `\d+` 匹配；如果是十六进制的，则用 `x[0-9a-fA-F]+` 匹配，最后的 `;` 也是必须出现的。注意这是一个捕获型括号，匹配完成后通过编号为 1 的分组可以直接提取出编码值。完整代码见例 9-10。

#### 例 9-10 使用 replacement 函数转换 HTML Entity

```
def htmlEntity2Char(matchObj) :
    if matchObj.group(0).startswith("&#x"):
```

```
#十六进制
return unichr(int(matchObj.group(1)[1:], 16))
#上一行的[1:]是为了去掉开头的 x
else :
    #十进制
    return unichr(int(matchObj.group(1)))

print re.sub(r"&#x?([0-9a-fA-F]+);", htmlEntity2Char, "&#x6536;&#x53D1;")
收发

Print.re.sub(r"&#x?([0-9a-fA-F]+);", htmlEntity2Char, "&#25910;&#21457;")
收发
```

在正则表达式替换时，将 *replacement* 参数设定为回调函数的做法，Python 和 JavaScript 都支持，.NET 也支持，不过.NET 通过 *deletgate*（可以理解为函数指针）实现回调函数的功能。具体的做法请参考第 10 章。

请注意，在实际开发中，应当谨慎采用替换中使用回调函数的做法，因为这样回调函数有能力执行更复杂的代码，所以可能会给恶意代码留下漏洞，带来安全性问题。

9.2.4 切分

切分也是正则表达式的常见操作之一，切分操作一般需要两个参数，即 *regex* 和 *string*，它以 *regex* 能匹配的文本为间隔，将 *string* 切分开来。它返回的通常是一个数组，元素是切分之后的片段。表 9-8 中列出了常用语言中的切分方法。

表 9-8 常用语言中的切分方法

语言	方法	备注
.NET	<b>Regex.Split</b> ( <i>input</i> , <i>pattern</i> )	静态方法
	<b>Regex.Split</b> ( <i>input</i> )	<b>Regex</b> 实例事先生成
Java	<b>String.split</b> ( <i>regex</i> )	静态方法
	<b>Pattern.split</b> ( <i>input</i> )	<b>Pattern</b> 实例事先生成
JavaScript	<b>string.split</b> ( <i>RegExp</i> )	
PHP	<b>preg_split</b> ( <i>pattern</i> , <i>subject</i> )	
Python	<b>re.split</b> ( <i>pattern</i> , <i>string</i> )	
Ruby	<b>string.split</b> ( <i>Regexp</i> )	

切分操作的简单示例是切分单词，尤其是英文文本的单词。英文文本用空格分隔各个单词，比如 **This is a example of using regex**，其中包含 7 个单词，用 6 段空白分隔——其中的空白可能是空格符，也可能是制表符，还有可能是换行符。不过，无论这些空白字符是什么，都可以由正则表达式 **\s+** 匹配；所以以这些空白为间隔，就可以把文本分隔为 7 段，每段一个单词，如例 9-11 所示。

**例 9-11 正则表达式切分**

```
print re.split(r"\\s+", "This is a example of using regex")
['This', 'is', 'a', 'example', 'of', 'using', 'regex']
```

这只是最简单的情况，似乎不用正则表达式也不难做到，其实依据正则表达式的切分其实可以完成更复杂的任务，比如将一个中英文混排的句子按照标点符号切分成多个分句。中英文混排文本如下：

I have read Tao, and learned a lot from it. 我读了《道德经》，从中学到了很多。

希望将它切分为 4 个分句：

```
I have read Tao
and learned a lot from it
我读了《道德经》
从中学到了很多
```

如果仅仅使用字符串操作，这几乎是不可能完成的，因为既要知道识别全角、半角的逗号和句号用来切分，还要知道书名号虽然是标点符号，但不能用作切分。但使用正则表达式则非常简单，因为正则表达式提供了对 Unicode Property 的支持（因为 Python 不支持 Unicode Property，所以这里使用 Java 举例），Unicode Property 按照字符的功能分类，`\p{Po}` 表示“除横线、括号、引号和连接符之外的任何标点符号”（☞ 123），代码见例 9-12。

**例 9-12 借助 Unicode Property 切分**

```
String str = "I have read Tao, and learned a lot from it. 我读了《道德经》，从中学到了很多。";
for (String s : str.split("\\p{Po}")) {
    System.out.println(s);
}
```

```
I have read Tao
and learned a lot from it
我读了《道德经》
从中学到了很多
```

最后补充简要介绍切分的次数限制。通常，切分会在所有匹配发生的地方进行，假设正则表达式能匹配  $n$  次，就切分  $n$  次，结果数组包含  $n+1$  个元素。不过，通常也可以用一个可选参数限制切分，将它设定为一个小于  $n$  的正数，则会进行  $n-1$  次切分（只有 Python 是例外，它会切分  $n$  次），返回数组的最后元素包含了“正则表达式第  $n-1$  次匹配右侧的所有文本”。

有的时候必须要用到这种功能，在例 9-13 中，SQL 日志文件的每一行都包含三个字段，分别是发生时间、执行时间、SQL 语句，以逗号分隔，而 SQL 语句中可能就包含逗号，这时需要显式限定切分次数，否则就会麻烦很多。

例 9-13 借助 Unicode Property 切分

```
logline = "14:59:35,0.133,select name,title from employee"
#没有设定切分次数
print re.split(r",", logline)
['14:59:35', '0.133', ' select name', 'title from employee']
#显式设定切分次数
print re.split(r",", logline, 2)
['14:59:35', '0.133', ' select name,title from employee']
```

## 9.3 正则表达式的优化建议

许多人认为正则表达式很强大、很快，现在的电脑运行速度很快，所以效率是无所谓的，但不同的正则表达式，不同的使用方法，所消耗的时间/空间可能相差几十倍甚至上百倍，在某些环境下（比如在移动设备上），这非常重要。本节整理了一些关于正则表达式的优化建议，合理使用它们有利于提高效率<sup>1</sup>。

在阅读和掌握这些建议之前，应当明确两点：首先，在正则表达式的匹配中，正确性是优先于效率的，常见的步骤是先保证正确，再优化效率；其次，如果效率不是严重的问题，就应当仔细审视优化的价值，因为有些优化技巧是需要改动正则表达式，影响可读性的，这么做是否值得，取决于具体的情况。

### 9.3.1 使用缓存

第 8 章已经介绍过正则表达式的匹配原理，生成正则表达式对象（也就是从字符串形式的表达式生成自动机）是非常复杂的过程。所以如果一个表达式需要反复使用，最好的办法就是将生成的表达式对象缓存起来。

在函数式处理中，正则表达式的缓存一般由语言自动完成的，比如 PHP 就可以缓存 4000 多个表达式。不过，这种机制对用户是透明的，一般用户完全不必关心。

在面向对象式处理中，正则表达式对象的生成可以由用户手工控制，例如 Python 中的 `re.compile(pattern)`、Java 中的 `Pattern.compile(regex)`、.NET 中的 `new Regex(regex)`。如果正则表达式对象需要反复用到，就应当用变量其保存起来，每次直接使用，避免重复编译。

---

<sup>1</sup> 为了节省篇幅，这里没有给出真正程序的测试数据，有兴趣的读者可以参考《精通正则表达式》（第 3 版）第 6 章的“性能测试”。

比如 Java 中的正则表达式替换，最简单的办法是 `String.replaceAll(regex, replacement)`，这是 `String` 类的静态方法，它的参数 `regex` 和 `replacement` 都是 `String` 对象。但是另一方面，也可以调用 `Matcher.replaceAll(replacement)` 方法，`Matcher` 对象由编译好的 `Pattern` 对象和 `String` 对象生成。如果要对多个 `String` 进行替换，先编译好 `Pattern` 对象，再对每个 `String` 生成 `Matcher`，调用成员方法 `Matcher.replaceAll(replacement)`，效率就高很多。

而 .NET 中还提供了更精确的缓存控制，可以人工设定 `Regex.CacheSize` 的值，指定缓存最近的  $n$  个正则表达式，根据场景选择合适的值，有可能大大提高效率。

### 9.3.2 尽量准确地表达意图

正则表达式中提供了许多结构（字符组、多选结构、量词、环视结构等等），许多结构的功能存在重叠，有些情况下，一些结构的功能可以完全由其他结构提供。这时候，不应该图省事使用“功能更全面”的结构，而应该恰当选取最合适（也就是更“专用”）的结构。

最常见的例子是多选结构和字符组，多选结构的功能完全覆盖了字符组。`(a|b|c|d)` 和 `[abcd]` 能匹配的文本是完全一样的，但是，多选结构是用于处理“可能出现若干个表达式”的，而字符组是专用于处理“可能出现若干个字符”的。针对例子中的这种情况，字符组的效率会远远高于多选结构，所以应当使用 `[abcd]` 而不是 `(a|b|c|d)`。

当然，字符组也不是任何时候都推荐使用。有些人不明白转义的规则，但知道字符组内部除去-字符，所有元字符都失去了特殊含义，就用字符组来“转义”，比如要匹配点号`.`，不使用`\.`，而是用`[.]`，这种做法同样会影响效率。

在第3章讲过，括号一旦使用，无论是用于分组，还是用于多选结构，都会产生副作用——匹配时会把文本保存起来，匹配完成之后通过对应的编号访问。但是许多时候，括号仅仅用来分组，匹配完成之后并不关心括号内的子表达式匹配的文本。这时候，就应当把捕获型括号`(regex)`改为非捕获型括号`(?:regex)`。不过，捕获型括号`(regex)`更美观，也更符合习惯，所以我推荐的做法是，在编写正则表达式时，先使用捕获型括号，确保准确无误之后，再改为非捕获型括号。

另外，如果要匹配任意的英文大小写字母，既可以使用`[a-zA-Z]`，也可以使用`[a-z]`，然后指定不区分大小写模式。相对来说，前者更明确，所以效率也更高些（当然这个例子最合适的是匹配“任意”字母的情况，如果要匹配的是确定单词，比如 `cat`，明确写出就是`[cC][aA][tT]`，明显不如 `cat` 直观，更不用说 `tomorrow` 之类的单词了）。但是，只有在要处理的文本特别多的情况下，这种优化才会带来比较大的效率提升，所以这个技巧应当斟酌使用，根据具体情况进行选择。

### 9.3.3 避免重复匹配

正则表达式非常灵活，同一个表达式往往可以匹配很多种不同形态的文本。但是，如果一个

表达式中包含若干个“能匹配不同形态文本”的子表达式，而这些子表达式能匹配的文本又有重叠，就可能造成效率问题。

最简单的例子是表达式 `a+a+`。上一章讲解过正则表达式的匹配原理：因为 `a+` 中 `+` 的存在，匹配时每遇到一个字符 `a` 都需要保存备用状态，以应付可能的回溯。这样，对于同一个字符，两个 `a+` 保存的备用状态数量就比单个 `a+` 要多出很多，而这其实是不必要的（某段文本如果不能被第一个 `a+` 匹配，必定也不能被第二个 `a+` 匹配），所以应当避免出现这样的表达式。

这类问题看起来很简单，在正常情况下很难有人写 `a+a+` 之类的表达式，甚至 `a+a*` 也不会出现，但是 `[0-9]+\w+` 之类的表达式则比较容易遇到，其中 `[0-9]` 和 `\w` 存在重复匹配的问题。更复杂的情况是类似 `(...[0-9]+)(\w+...)` 之类的表达式，看起来 `(...[0-9]+)` 和 `(\w+...)` 的逻辑都足够清楚，各部分前后连接的地方却出现了重复匹配的问题。在构造复杂表达式的过程中，很容易出现这类问题，一定要多加注意。

更复杂的重复匹配的问题来自多选结构。之前提到过，在多选结构中一定要尽力避免各个多选分支可能匹配相同的文本。在之前提到的匹配双引号字符串的例子中，使用多选结构 `(\\.|[^\\""])` 匹配起始双引号和结束双引号之间的文本，其中第一个多选分支 `\\.` 匹配任何转义序列，第二个多选分支本意是匹配其他字符（也就是除去 `"` 之外的其他字符），不仔细思考可能会使用字符组 `[^"]`，这样不但会降低效率，还会造成匹配错误。

重复匹配的情况看似麻烦，其实并非如此。总结起来，主要有两条规则，如果能做到，则大部分重复匹配的问题都可以避免。

(1) 凡类似 `regex1*regex2*`（量词不限于 `*`）的表达式，都要尽量避免 `regex1` 和 `regex2` 能匹配相同的内容；

(2) 凡类似 `(regex1|regex2)` 的表达式，都要尽量避免 `regex1` 和 `regex2` 能匹配相同的内容。

### 9.3.4 独立出文本和锚点

正则表达式包含多种复杂的结构，提供了匹配文本的各种可能性。可能性增加的代价就是降低了处理的速度。根据之前讲解的知识，正则表达式匹配都要经过“编译正则表达式”和“逐次试错”的过程，普通的字符串操作没有这么复杂，按照目前普遍采用的 Boyer-Moore 字符串检索算法<sup>1</sup>，能很快地找到某个字符串在文本中的位置。

对正则表达式来说，我们常常可以发现多种“可能性”的共同点，也就是各种可能中都具有的，不变的字符串，将它们提取出来，就可以依靠普通的字符串操作，迅速定位到某个位置，再“逐次试错”。

<sup>1</sup> 请参考 [http://en.wikipedia.org/wiki/Boyer-Moore\\_string\\_search\\_algorithm](http://en.wikipedia.org/wiki/Boyer-Moore_string_search_algorithm)。

仔细观察正则表达式 `(this|that|these|those)` 可以发现，多选结构的所有多选分支都以 `th` 开头，所以这个表达式如果能够匹配，前两个字符必然是 `th`。但正则表达式处理程序未必能“意识”到这一点，所以不妨人工将 `th` 提取出来，将表达式改为 `th(is|at|ese|ose)`。同样，如果多选分支的结尾是相同的，依然可以人工提取出来，把表达式 `(option|action)` 改为 `(op|ac)tion`，也有这种效果<sup>1</sup>。

将文本独立出来的技巧也适用于量词。比如 `=+` 要求等号 `=` 至少必须出现一次，它等价于 `=*`，但是后者独立出了一个固定的字符，所以效率更高；再比如 `={4}`，它等价于 `====`，但后者独立出 4 个固定的字符，因此效率要更高。测试数据显示，在 Python 和 Java 中，后者的效率是前者的 100 倍以上；不过，Ruby 和 .NET 能自动识别两者，自动替你进行优化，所以两者效率是一样的。

将文本独立出来可以提高效率，同样道理，将锚点独立出来也可以提高匹配的效率。比如 `(^this|^that)`，正则表达式处理程序未必知道两个分支都必须在字符串的起始位置（如果指定了多行模式，则是行的起始位置）开始匹配。如果将 `^` 独立出来，改为 `^(this|that)`，也可以提高效率。

## 9.4 别过分依赖正则表达式

正则表达式的功能很强大，有些人学会了之后随时都想使用它。但是，过分依赖正则表达式并不是好的习惯。下面列举几个过分依赖正则表达式的常见错误。

### 9.4.1 彻底放弃字符串操作

最常见的错误是，某个任务明明只需要字符串操作就可以完成，但非要使用正则表达式。上一章已经分析过正则表达式的匹配原理，正则表达式匹配需要先构建有穷自动机，这种操作的成本高昂。如果只需要查找；同样，如果只需要判断字符串开头以某个固定子串开头，以某个固定子串结尾，一般语言中都提供了 `startswith()`、`endswith()` 之类的方法，相比之下，正则表达式即便使用了锚点，效率也要低得多。表 9-9 比较了这两类操作。

表 9-9 这些情况下，字符串操作比正则表达式更高效

字符串操作	正则表达式操作
<code>string.index("cat") != -1</code>	<code>re.search("cat", string) != None</code>
<code>string.startswith("cat")</code>	<code>re.search("^cat", string) != None</code>
<code>string.endswith("cat")</code>	<code>re.search("cat\$", string) != None</code>

如果检查多个固定形式的字符串，正则表达式的效率仍然比不上简单的字符串操作。比如要

<sup>1</sup> 这样的代价是阅读代码的难度增加了。所以我的建议是：这类优化措施，首先必须保证正确；其次，适当的时候应当添加注释。

检查一个字符串中是否包含 `one`、`two`、`three` 这三个单词中的任意一个，正则表达式的解法是查找 `(one|two|three)`，但它的效率比不上依次判断是否存在 `one`、`two`、`three`。因为简单字符串查找已经有很多高效的算法，“大而全”的正则表达式却必须在文本的每个位置尝试匹配。要查找的字符串越多，正则表达式的效率就越低。在实际应用中，经常可能要检查文本中是否包含某些不期望的文本（“敏感”词），如果文本很长，不期望出现的词很多，使用正则表达式就可能严重影响效率。

## 9.4.2 思维定势

另一个常见的错误是形成了思维定势，解决问题的思路被正则表达式的 API 限制。针对正则表达式的常见操作（提取、验证等），编程语言和工具都提供了对应的 API，但这并不是说只能硬性使用这些 API。

回过头看之前讲解的密码字符串匹配的例子，符合要求的密码必须满足五点要求：

- （1）密码的长度在 6~12 个字符之间；
- （2）只能由小写字母、阿拉伯数字、横线组成；
- （3）开头和结尾不能是横线；
- （4）不能全部是数字；
- （5）不容许有连续（2 个及以上）的横线。

单独针对其中每个条件判断，表达式都很容易写（甚至更容易，比如“不容许有连续（2 个及以上）的横线”，可以先验证表达式中是否存在连续的横线，再对结果取“非”）。总的验证如果要成立，则必须同时满足这 5 个条件，也就是对 5 个单独验证的结果进行“与”运算。这种思路的代码见例 9-14，比起费力把 5 个条件整合到一个表达式中，这种办法要清晰多了。

### 例 9-14 用一组正则表达式验证

```
cond1Regex = r"\A.{6,12}\Z"
cond2Regex = r"\A.[0-9A-Za-z-]*\Z"
cond3Regex = r"(\A-|- \Z)" #开头或结尾是横线
cond4Regex = r"\D" #出现一个非数字字符
cond5Regex = r"--" #连续的横线
#Python 中语句末尾没有;，行末的\表示同一条语句未结束
re.search(cond1Regex, password) != None \
and re.search(cond2Regex, password) != None \
and re.search(cond3Regex, password) == None \ #注意，此处取非
and re.search(cond4Regex, password) != None \
and re.search(cond5Regex, password) == None \ #注意，此处取非
```

另一个例子是去掉行首尾的空白字符。第 66 页已经说明，最简单的办法是分两步处理，在多行模式下，分别将 `^\s+` 和 `\s+$` 能匹配的文本替换为空字符串，这种做法直观，而且效率最高；



如果非要用一个表达式替换来完成，不但可能出错，效率也无法保障，甚至可能降低很多。

从这点推广开来，正则表达式并没有规定，执行某类操作就必须使用对应的 API，同一个任务有多种完成方法。比如从日期字符串 `2010-10-22` 提取出年、月、日信息，既可以用表达式 `(\d{4})-(\d{2})-(\d{2})` 匹配再提取对应捕获分组，也可以用 `_` 切分，两者的结果是一样的。

### 9.4.3 正则表达式可以匹配各种文本

最后一个常见错误是，以为正则表达式能完成所有的任务。正则表达式的功能虽然强大，但只限于描述文本特征，而且甚至无法覆盖所有的文本特征。

如果要匹配这样的文本：`01` 或 `0011` 或 `000111`……总之是由 `0` 和 `1` 构成的字符串，连续的 `0` 在前，连续的 `1` 在后，而且 `0` 的数目和 `1` 的数目应当相等。在正则表达式中，对于这类情况应当使用量词，所以表达式应当类似 `0+1+`。遗憾的是，在正则表达式中我们无法表示“限定 `0` 的量词等于限定 `1` 的量词”，所以仅仅依靠正则表达式是无法解决这类问题的。而例 9-15 的逻辑则要简单很多：先用正则表达式 `(0+)(1+)` 验证，通过之后，再对比这个表达式中编号为 1、2 的两个捕获分组匹配文本的长度，如果相等，则留下，否则忽略。

例 9-15 正则表达式配合其他运算进行验证

```
def validate(str) :
    matchObj = re.search(r"\A(0+)(1+)\Z", str)
    return (matchObj != None) and (len(matchObj.group(1)) == len(matchObj.group(2)))

validate("001")    # => False
validate("0011")   # => True
```

还有些时候，要匹配的文本不只具有文本特征，还有其他方面的意义要求，这时候正则表达式更是无能为力了，闰年的匹配就是如此。

一般来说，闰年就是可以被 4 整除的年份。闰年的两位尾数其实很好判断：如果十位是 0、2、4、6、8，则个位可以是 0、4、8；如果十位是 1、3、5、7、9，则个位可以是 2、6，看来，匹配闰年的表达式是 `\d{2}([02468][048]|[13579][26])`（这里暂时不考虑两端的断言）。

但是，闰年还有另一条规定：尾数为 00 的年份，除非整个数字能被 400 整除的年份，否则不是闰年。所以 2000 是闰年，而 1900、1800 都不是闰年。如果要提取某段文本中所有闰年的年份，单纯使用正则表达式就非常吃力。但如果先提取出所有“疑似闰年”的年份，再用一个函数判断这个年份如果尾数是 00 能否整除 400，对之前提取的年份进行过滤，则要简单许多，代码见例 9-16。

## 例 9-16 正则表达式配合其他运算验证闰年

```
def validateLeapYear(str) :  
    return (re.search(r"\A\d{2}([02468][048]|[13579][26])\Z", str) != None) and  
int(str) % 400 != 0  
  
validateLeapYear("1986")    # => False  
validateLeapYear("1988")    # => True  
validateLeapYear("2000")    # => False  
validateLeapYear("2100")    # => True
```



# 第三部分

## 第10章 .NET

在微软的.NET Framework 里的 C#、C++、VB.NET 等各种语言中，正则表达式的包是统一的，都是 `System.Text.RegularExpressions`；而且，基本的对象和方法在各语言中也没有太多差异，所以在某种语言中编写的关于正则表达式的程序，可以很方便地移植到.NET 的其他语言中。为节省篇幅，本章的大多数示例使用 C#，VB.NET 的代码与之没有本质差别，所以不再重复<sup>1</sup>。

### 10.1 预备知识

.NET 语言中的正则表达式相关类都存在于 `System.Text.RegularExpressions` 中，如果要使用正则表达式，必须首先导入它，最简单的办法是在程序文件头部这样写：

```
C#
using System.Text.RegularExpressions;

VB.NET
Imports System.Text.RegularExpressions;
```

对应正则表达式的类名称为 `Regex`，这个类会在下面详细介绍，这里先介绍.NET 中的 `Regex.IsMatch(input, regex)` 方法<sup>2</sup>，它返回一个布尔值，表示 `regex` 对应的表达式能否在 `input` 中找到匹配。即使只能匹配一个子串，也会返回 `True`，如果要检查整个字符串能否由表达式匹配，可以在表达式首尾加上 `\A` 和 `\z`。

```
Regex.IsMatch("1", "\\d");           // => True
Regex.IsMatch("1a", "\\d");          // => True
Regex.IsMatch("1", "\\A\\d\\z");      // => True
Regex.IsMatch("1a", "\\A\\d\\z");    // => False
```

请注意其中的字符组简记法和锚点都使用了两个反斜线字符，因为需要转义。其实 C# 的 `String` 也提供了原生字符串（Verbatim String），它不需要考虑字符串转义，所以非常适合表示正则表达式：正则表达式是怎样的，原生字符串就写成怎样。具体做法是在 `String` 开头的双引号之前添加 `@`，为方便阅读，本章的示例代码都采用这种写法。

```
Regex.IsMatch("1", @"\\d");           // => True
Regex.IsMatch("1a", @"\\d");          // => True
Regex.IsMatch("1", @"\A\d\z");        // => True
```

<sup>1</sup> 详细文档可参考 <http://msdn.microsoft.com/en-us/library/hs600312%28v=vs.71%29.aspx>。

<sup>2</sup> .NET 中的正则 API 的参数中，字符串在前，正则表达式在后，顺序不同于 Python、Java、PHP 等语言，如果要同时在多种语言中编程，请留意这一点。

```
Regex.IsMatch("1a", @"\\A\\d\\z"); // => False
```

因为原生字符串不识别任何反斜线转义序列，字符串中使用的\\n、\\t 表示法都无法用在原生字符串中，所以本章只对正则表达式参数使用原生字符串，而不对字符串参数使用原生字符串。

VB.NET 中的字符串比较特别，处理字符串文字时不识别任何转义序列，所以可以理解为，VB.NET 中的所有字符串天生就是原生字符串，唯一例外的文字是紧密相连的两个双引号""，它表示一个双引号"。下面给出了几个 VB.NET 中的例子。

```
Regex.IsMatch("1", "\\d"); // => True
Regex.IsMatch("""", """""); // => True
```

## 10.2 正则功能详解

### 10.2.1 列表

表 10-1 中列出了.NET 中的正则功能。

表 10-1 .NET 中的正则功能列表

功能	记法	说明
字符组 ④1	[...] [^...]	完全支持 Unicode
POSIX 字符组 ④15	[[:digit:]]	不支持
Unicode 属性 ④121	\p{...}	完全支持 Unicode
字符组简记法 ④115	\d \D \w \W \s \S	采用 Unicode 规则
行起始位置 ④61	^ \A	支持
行结束位置 ④61	\$ \z \Z	支持
单词边界 ④117	\b \B	采用 Unicode 规则
顺序环视 ④68	(?=...) (?!...)	支持
逆序环视 ④68	(?<=...) (?<!...)	支持
匹配模式 ④81	i m s x	支持
模式作用范围 ④89	(?modifier) (?-modifier)	支持
纯文本模式 ④99	\Q...\E	不支持
捕获分组及引用④44	(...) \num \$num	支持
命名分组 ④53	(?<name>...) \k{name} \${name}	支持
非捕获分组 ④54	(?:...)	支持
多选结构 ④39	(... ...)	支持
匹配优先量词 ④18	? * +	支持
忽略优先量词 ④26	?* ?+ {n,m}?	支持

## 10.2.2 字符组

在.NET 语言中,所有的字符都是 Unicode 字符,所以在字符组中可以直接使用中文,因而不会出现多字节字符错误匹配的问题;另一方面,因为所有字符都是 Unicode 字符,所以点号也可以匹配中文字符。

```
Regex.IsMatch(".", @"A 你\z"); // => True
Regex.IsMatch("罚", @"A[收发]\z"); // => False
```

值得注意的是, .NET 语言中的字符组可以进行集合运算,此功能非常适合进行字符组的“减”运算——如果需要匹配英文中所有小写的辅音字母,也就是从 26 个小写字母中“减去”5 个元音字母,可以写作 `[a-z-[aeiou]]` (注意不是 `[[a-z]-[aeiou]]`)。它的意思是,“从 a 到 z 的所有字符(也就是小写英文字符)”与“除 aeiou 之外的任何字符”取交集,实际结果就是所有的小写辅音字母。

```
Regex.IsMatch("a", @"[a-z-[aeiou]]"); // => False
Regex.IsMatch("b", @"[a-z-[aeiou]]"); // => True
```

也可以用这个字符组来匹配“英文大小写字母和数字(不包括下划线)”。

```
Regex.IsMatch("a", @"[\w-[_]]"); // => True
Regex.IsMatch("0", @"[\w-[_]]"); // => True
Regex.IsMatch("_", @"[\w-[_]]"); // => False
```

.NET 中可以用 `\uhex` 指定 Unicode 码值,其中 *hex* 是 4 位十六进制数。所以,在 Unicode 编码环境下可以用字符组 `[\u4e00-\u9fff]` 匹配所有中文字符。

```
Regex.IsMatch("我", @"[\u4e00-\u9fff]"); // => True
```

## 10.2.3 Unicode 属性

.NET 对 Unicode 字符组的支持比较好,它支持 Unicode Property。这样,用 `\p{N}` 就可以匹配所有的数字字符,包括中文的全角数字字符,比如 0、1、2…; `\p{P}` 也可以匹配各种标点符号,包括中文的冒号、书名号《和》等。

```
Regex.IsMatch("1", @"\p{N}"); // => True
Regex.IsMatch(":", @"\p{P}"); // => True
Regex.IsMatch("《", @"\p{P}"); // => True
```

.NET 也支持 Unicode Block,其记法是以 Is 为前缀的,这与 Java 不一样,使用时应注意<sup>1</sup>,比如匹配中文字符的 Unicode Block,就应当写作 `IsCJKUnifiedIdeographs`。

```
Regex.IsMatch("我", @"\p{IsCJKUnifiedIdeographs}"); // => True
```

<sup>1</sup> Unicode Block 的细节可参考 <http://msdn.microsoft.com/en-us/library/20bw873z.aspx#SupportedNamedBlocks>。

关于 Unicode 属性的细节，请参考第 123 页。

## 10.2.4 字符组简记法

在 .NET 的正则表达式中，`\d`、`\D`、`\w`、`\W`、`\s`、`\S` 都使用 Unicode 匹配规则。也就是说，`\d` 不仅仅等价于 `[0-9]`，还可以匹配其他语言中的数字字符；`\w` 不仅仅等价于 `[0-9a-zA-Z_]`，还可以匹配其他语言中的文字字符；`\s` 也可以匹配 ASCII 编码之外的空白字符。

```
//半角字符
Regex.IsMatch("1", @"\d"); // => True
Regex.IsMatch("a", @"\w"); // => True
Regex.IsMatch(" ", @"\s"); // => True

//全角字符及中文字符
Regex.IsMatch(" 1", @"\d"); // => True
Regex.IsMatch("我", @"\w"); // => True
Regex.IsMatch(" ", @"\s"); // => True
```

如果要恢复这几个字符组简记法的 ASCII 匹配规则，可以指定使用 ECMAScript 模式，因为它没有对应的模式修饰符，所以只能以预定义常量 `RegexOptions.ECMAScript` 的方式指定。

```
//全角字符及中文字符
Regex.IsMatch(" 1", @"\d", RegexOptions.ECMAScript); // => False
Regex.IsMatch("我", @"\w", RegexOptions.ECMAScript); // => False
Regex.IsMatch(" ", @"\s", RegexOptions.ECMAScript); // => False
```

## 10.2.5 单词边界

.NET 正则表达式中的单词边界采用 Unicode 匹配规则，也就是说，如果 `\b` 能匹配，则一侧必须出现 `\w` 能匹配的字符，另一侧不能出现 `\w` 能匹配的字符（不是“另一侧出现 `\w` 不能匹配的字符”）。

```
//英文字符-半角标点
Regex.IsMatch("a,", @"a\b"); // => True
//英文字符-全角标点
Regex.IsMatch("a ", @"a\b"); // => True
//英文字符-结束
Regex.IsMatch("a", @"a\b"); // => True
//中文字符-半角标点
Regex.IsMatch("我,", @"我\b"); // => True
//中文字符-全角标点
Regex.IsMatch("我 ", @"我\b"); // => True
//中文字符-空字符串
Regex.IsMatch("我\b ", @"我"); // => True
//全角标点-结束（未出现单词字符）
```



```
Regex.IsMatch(" ", "@", "\b"); // => False
//中文字符-英文字符（未出现单词字符）
Regex.IsMatch("我 a", "@我\b"); // => False
```

指定 ECMAScript 模式会影响 `\w` 的匹配，所以 `\b` 的匹配也同样会受到影响。

```
//英文字符-半角标点
Regex.IsMatch("a", "@a\b", RegexOptions.ECMAScript); // => True
//英文字符-全角标点
Regex.IsMatch("a ", "@a\b", RegexOptions.ECMAScript); // => True
//英文字符-结束
Regex.IsMatch("a", "@a\b", RegexOptions.ECMAScript); // => True
//中文字符-半角标点
Regex.IsMatch("我", "@我\b", RegexOptions.ECMAScript); // => False
//中文字符-全角标点
Regex.IsMatch("我 ", "@我\b", RegexOptions.ECMAScript); // => False
//中文字符-空字符串
Regex.IsMatch("我\b ", "@我", RegexOptions.ECMAScript); // => False
//全角标点-结束（未出现单词字符）
Regex.IsMatch(" ", "@", "\b", RegexOptions.ECMAScript); // => False
//中文字符-英文字符（未出现单词字符）
Regex.IsMatch("我 a", "@我\b", RegexOptions.ECMAScript); // => True
```

## 10.2.6 行起始/结束位置

`^`匹配的是“行开始的位置”，默认情况下它只能匹配“整个字符串的开始位置”，如果指定使用多行模式（Multiline Mode），`^`可以匹配字符串内部的文本行的开始位置；而 `\A` 无论在什么情况下，都匹配“整个字符串的开始位置”。

```
Regex.IsMatch("1\n2\n", @"^1"); // => True
Regex.IsMatch("1\n2\n", @"^2"); // => False
Regex.IsMatch("1\n2\n", @"(?m)^2"); // => True

Regex.IsMatch("1\n2\n", @"\A2"); // => True
Regex.IsMatch("1\n2\n", @"\A2"); // => False
Regex.IsMatch("1\n2\n", @"(?m)\A2"); // => False
```

在默认情况下，.NET 中的 `$`、`\Z`、`\z` 匹配的是整个字符串的结束位置（如果结束位置有换行符，则 `$` 和 `\z` 匹配这个换行符之前的位置），多行模式只会影响到 `$` 的匹配，在这种模式下它可以匹配文本内部行的结束位置。所以如果要做严格准确的验证，应当使用 `\z`。

```
Regex.IsMatch("1\n", @"1$"); // => True
Regex.IsMatch("1\n2", @"1$"); // => False
Regex.IsMatch("1\n2", @"(?m)1$"); // => True

Regex.IsMatch("1\n", @"1\Z"); // => True
```

```
Regex.IsMatch("1\n2", @"1\Z");           // => False
Regex.IsMatch("1\n2", @"(?m)1\Z");        // => False

Regex.IsMatch("1", @"1\z");               // => True
Regex.IsMatch("1\n", @"1\z");             // => False
Regex.IsMatch("1\n", @"(?m)1\z");         // => False
```

10.2.7 环视

.NET 中的肯定顺序环视 `(?=...)` 和否定顺序环视 `(?!...)` 内可以使用任意形式的子表达式。

```
Regex.IsMatch("cab", @"c(?:=ab+)");       // => True
Regex.IsMatch("ccd", @"c(?:=(ab+|cd))");   // => True
Regex.IsMatch("cabbbb", @"c(?:=(ab+|cd))"); // => True
```

.NET 的逆序环视值得大加表扬——在目前所有的编程语言中，只有.NET 对逆序环视没有任何限制，也就是说，可以在.NET 的逆序环视中使用任何形式的表达式。

```
Regex.IsMatch("ab", @"ab(?:<=ab)");       // => True
Regex.IsMatch("cd", @"cd(?:<=((ab)*cd|cd))"); // => True
Regex.IsMatch("cd", @"cd(?:<=c+d?)");     // => True
```

虽然.NET 对逆序环视的支持相当完备，但我并不推荐在逆序环视中使用特别复杂的表达式，因为这样可能会严重影响性能。

10.2.8 匹配模式

表 10-2 中列出了.NET 中的匹配模式。

表 10-2 .NET 中的匹配模式

常量	修饰符	说明
IgnoreCase	i	不区分 ASCII 字符的大小写
IgnorePatternWhitespace	x	容许正则表达式中出现注释，表达式中的空白字符，以及#开始到行末的文本，都视为注释
常量	修饰符	说明
Multiline	m	容许 <code>^</code> 和 <code>\$</code> 不仅仅匹配字符串的起始和结束位置，还可以匹配字符串内部文本行的起始和结束位置
Singleline	s	容许点号 <code>.</code> 匹配任何字符，包括换行符
ExplicitCapture	n	只为命名分组保留匹配的文本，非命名分组不保留匹配的文本
ECMAScript	无	指定字符组简记法 <code>\d</code> <code>\D</code> <code>\w</code> <code>\W</code> <code>\s</code> <code>\S</code> 采用 ASCII 匹配规则

( 续表 )

常量	修饰符	说明
Compiled	无	将正则表达式直接编译为底层 MSIL ( Microsoft Intermediate Language ), 会降低启动速度, 增加内存消耗, 但匹配速度能大幅提升, 适合将特别常用的正则表达式编译到 DLL 中重用

匹配模式可以在表达式中以 `(?modifier)` 指定, 也可以在生成 `Regex` 对象时, 将预定义常量作为参数指定。

```
Regex.IsMatch("A", @"a");           // => False
Regex.IsMatch("A", @"(?i)a");       // => True
Regex.IsMatch("A", @"a", RegexOptions.IgnoreCase); // => True
```

.NET 也支持用 `(?-modifier)` 停用某个模式。

```
Regex.IsMatch("aBb", @"a(?i)b(?-i)b"); // => True
Regex.IsMatch("aBB", @"a(?i)b(?-i)b"); // => False
```

10.2.9 捕获分组的引用

在.NET 中, 如果要在正则表达式内部引用捕获分组, 应当使用 `\num` 记法, 其中 `num` 为对应捕获分组的编号。

```
Regex.IsMatch("ab", @"([a-z])\1"); // => False
Regex.IsMatch("aa", @"([a-z])\1"); // => True
```

当然也可以使用命名分组, 这样看起来更清楚, 在表达式中引用命名分组的记法是 `k<name>`。

```
Regex.IsMatch("ab", @"(?<char>[a-z])\k<char>"); // => False
Regex.IsMatch("aa", @"(?<char>[a-z])\k<char>"); // => True
```

如果要在替换时引用捕获分组, 应当使用 `$num` 记法。

```
Regex.Replace("2010-12-20", @"(\d{4})-(\d{2})-(\d{2})", "$2/$3/$1");
12/20/2010
```

当然也可以使用命名分组, 但是在 `replacement` 字符串中引用的记法不同于在表达式中, 应当写为 `${name}`。

```
Regex.Replace ("2010-12-20", @"(?<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})", @"${month}/${day}/${year}");
12/20/2010
```

如果要在 `replacement` 字符串中使用 `$` 字符, 必须将它转义为 `$$`。

```
Regex.Replace("the price is 12.99", @"\d+\.\d{0,2}", "$$$0");
the price is $12.99
```

## 10.3 正则 API 简介

.NET 语言中的正则表达式相关类都存在于 `System.Text.RegularExpressions` 中，通常用到的是这几个类：`Regex`、`Match`、`Group`，关系如图 10-1 所示。`Regex` 对应正则表达式，一个 `Regex` 对象与一个 `String` 对象关联，生成一个 `Match` 对象，它表示 `Regex` 在 `String` 中的一次匹配，对它调用 `NextMatch()` 方法，可以获得下一次匹配的 `Match` 对象；借助 `Match` 获得对应的 `Group` 对象，可以访问当次匹配时的分组捕获信息。

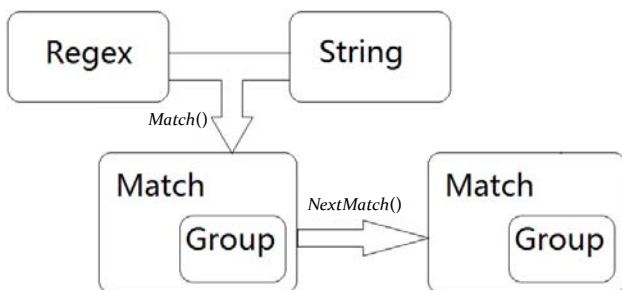


图 10-1 .NET 中与正则表达式相关的类

```

Regex regex = new Regex(@"\d{4}-\d{2}-\d{2}");
Match match = regex.Match("2010-12-20 2011-02-14");
int[] groupNums = regex.GetGroupNumbers();
while (match.Success) {
    Console.WriteLine(match.Groups[0].Value);
    match = match.NextMatch();
}
2010-12-20
2011-02-14

```

### 10.3.1 Regex

`Regex` 是 .NET 中的“正则表达式对象”（之前介绍过，虽然正则表达式是用字符串形式表示的，但它并不等于字符串）。要使用正则表达式，首先必须从字符串“编译”（也可以说“生成”）出 `Regex` 对象，常用方法是这样的：

```
Regex regex = new Regex(@"ab+");
```

如果要指定匹配模式，可以在表达式中使用修饰符 `(?modifier)` 指定，也可以使用预定义常量。下面的两个 `Regex` 对象虽然生成方式不同，却是等价的。

```

Regex regex = new Regex(@"(?i)ab+");
Regex regex = new Regex(@"ab+", Pattern.CASE_INSENSITIVE);

```

如果要同时指定多种模式，可以连写模式修饰符，也可以直接用|将预定义常量连接起来。下面的两个 `Regex` 对象也是等价的。

```
Regex regex = new Regex (@"(?is)ab+");
Regex regex = new Regex (@"ab+", RegexOptions.IgnoreCase | RegexOptions.Singleline);
```

下面介绍 `Regex` 类的主要成员方法。

### 10.3.1.1 Static String Escape(String text)

这个方法用来取消字符串 `text` 中所有转义字符的特殊含义。如果需要取消某个字符串中所有元字符的特殊含义（比如从外界读入的字符串，用于正则表达式搜索），就可以使用这个方法。

```
Regex.IsMatch("aacb", @"a*.b"); // => True
Regex.IsMatch("aacb", Regex.Escape(@"a*.b")); // => False

Regex.IsMatch("a*.b", @"a*.b"); // => False
Regex.IsMatch("a*.b", Regex.Escape(@"a*.b")); // => False
```

### 10.3.1.2 String[] GetGroupNames()

这个方法用来返回当前 `Regex` 对象中各命名分组的名字（对应整个表达式的默认分组名字是 0），如果只有捕获分组没有命名分组，则将捕获分组的数字编号 1、2、3 之类作为字符串返回。

```
Regex regex = new Regex(@"(?<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})");
foreach (String name in regex.GetGroupNames()) {
    Console.WriteLine(name + " ");
}

0 year month day

Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
foreach (String name in regex.GetGroupNames()) {
    Console.WriteLine(name + " ");
}

0 1 2 3
```

### 10.3.1.3 int[] GetGroupNumbers()

这个方法用来返回当前 `Regex` 对象中各捕获分组的编号，使用数字编号（`int` 类型变量）访问捕获分组的速度比使用分组名（`String` 类型变量）要快。

```
Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
foreach (int num in regex.GetGroupNumbers()) {
    Console.WriteLine(num + " ");
}

0 1 2 3
```

### 10.3.1.4 bool IsMatch(String input)

这个方法用来判断当前的 `Regex` 对象能否在字符串 `input` 中匹配某个子串。如果能则返回 `true`，否则返回 `false`。因为它并不检查整个字符串能否由正则表达式匹配，所以在数据验证，必须在表达式首尾加上锚点 `\A` 和 `\z`。

```
(new Regex(@"\d")).IsMatch("a1");           // => True
(new Regex(@"\d")).IsMatch("ab");           // => False

(new Regex(@"\A\d\z")).IsMatch("a1");       // => False
(new Regex(@"\A\d\z")).IsMatch("1");       // => True
```

它有两个相应的静态便捷函数：

```
Regex.IsMatch(String input, String pattern);
Regex.IsMatch(String input, String pattern, RegexOptions options);
```

### 10.3.1.5 String Replace(String input, String replacement)

这个方法用来将字符串 `input` 中 `Regex` 对象能替换的字符串都替换为 `replacement`。

```
Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
Console.WriteLine(regex.Replace("2010-12-20 2011-02-14", "$2/$3/$1"));
12/20/2010 02/14/2011
```

在默认状态下会替换对所有能匹配的子串；也可以通过设定第三个参数 `count`，人工指定替换可能发生的最大次数。

```
Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
Console.WriteLine(regex.Replace("2010-12-20 2011-02-14", "$2/$3/$1", 1));
12/20/2010 2011-02-14
```

它有三个对应的静态便捷函数：

```
Regex.Replace(String input, String pattern, String replacement);
Regex.Replace(String input, String pattern, String replacement, int count);
Regex.Replace(String input, String pattern, String replacement, RegexOptions option);
```

`Replace()` 方法也可以不指定 `replacement` 字符串，而是提供 `MatchEvaluator` 对象（可以把它理解为函数指针），它接收 `Match` 对象，返回字符串。这样，替换操作就可以更加灵活。下面的代码展示了如何将单词统一为首字母大写格式。

```
static string Capitalize(Match match)
{
    //正则表达式已经用两个捕获分组分开第一个字母和之后的字母
    return match.Groups[1].Value.ToUpper() + match.Groups[2].Value.ToLower();
}

Console.WriteLine(new Regex(@"(?i)\b([a-z])([a-z]+\b)").Replace("one TWO thREe",
new MatchEvaluator(Capitalize)));
One TWO THREE
```

10.3.1.6 String[] Split(String input)

这个方法用来切分字符串 *input*，以当前的 **Regex** 对象匹配的子串为分隔。

```
foreach(String s in (new Regex(@"\s+")).Split("one two three")) {
    Console.WriteLine(s);
}
```

one  
two  
three

在默认情况下，这个方法会切分尽可能多的次数；不过，也可以通过第二个参数 *count* 指定返回数组的大小。

```
foreach(String s in (new Regex(@"\s+")).Split("one two three", 2)) {
    Console.WriteLine(s);
}
```

one  
two three

表 10-3 总结了 *count* 取各种值的情况对结果的影响（未指定 *count* 时，返回包含 *n* 个元素的数组，切分 *n-1* 次）。

表 10-3 count 各种取值的意义

count 取值	意义
count < 0	不容许，报错
count = 0	与未设置 count 一样，切分 <i>n-1</i> 次
0 < count < <i>n</i>	返回包含 count 个元素的数组，切分 count-1 次，最后一个元素是第 count-1 次切分后右侧剩下的所有文本
count >= <i>n</i>	等于未指定 count

它有三个相应的静态便捷函数：

```
Regex.Split(String input, String pattern);
Regex.Split(String input, String pattern, int count);
Regex.Split(String input, String pattern, RegexOptions option);
```

关于匹配、替换、切分，**Regex** 对象都提供了对应的静态便捷函数，这样代码更加易懂，而且不用显式生成一堆对象；坏处则是每次调用都要重新检查 *pattern* 字符串，重新生成 **Regex** 对象。为解决这个问题，.NET 可以缓存最近的 *n* 个 **Regex** 对象，*n* 的默认值为 15，如果对此不满意，可以重新设置（将值设定为 0 表示禁用缓存）。

```
Regex.CacheSize = 20;
```

## 10.3.2 Match

**Match** 可以理解为“某次具体匹配的结果对象”，调用 **Match()** 方法，把编译好的 **Regex** 对象“应用”到某个 **String** 对象上，就获得了作为“本次匹配结果”的 **Match** 对象。

每个 **Match** 对象对应到一次具体的匹配，**Match** 对象的 **Success** 属性标识本次匹配是否成功，如果成功；还可以通过 **Index**、**Length**、**Value** 三个变量获得关于此次匹配的具体信息。

Index	本次匹配文本在字符串中的开始位置
Length	本次匹配文本的长度
Value	本次匹配文本的内容

如果想进行下一次匹配，可以通过 **nextMatch()** 方法获得新的 **Match** 对象。

```
Regex regex = new Regex(@"\d{4}-\d{2}-\d{2}");
Match match = regex.Match("2010-12-20 2011-02-14");
while(match.Success) {
    Console.WriteLine(match.Value);
    match = match.NextMatch();
}
```

2010-12-20

2011-02-14

如果正则表达式中包含捕获分组（包括普通捕获分组和命名分组），则返回的结果中会包含捕获分组匹配的信息，这些信息封装为一个 **GroupCollection** 对象，其中每个元素都是 **Group** 对象，它其实非常类似 **Match** 对象，也包括 **Success**、**Value**、**Length**、**Index** 等属性。如果需要某个捕获分组匹配的具体信息，可以直接访问对应的 **Group** 对象得到。

```
Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
Match match = regex.Match("2010-12-20 2011-02-14");
int[] groupNums = regex.GetGroupNumbers();
while (match.Success) {
    for (int i = 0, n = groupNums.Length; i < n; i++) {
        int groupNum = groupNums[i];
        Console.WriteLine(match.Groups[groupNum].Value + " starts at " +
            match.Groups[groupNum].Index + " and length is " +
            match.Groups[groupNum].Length);
    }
    match = match.NextMatch();
}
```

2010-12-20 starts at 0 and length is 10

2010 starts at 0 and length is 4

12 starts at 5 and length is 2



```

20 starts at 8 and length is 2
2011-02-14 starts at 11 and end at 21
2011 starts at 11 and length is 4
02 starts at 16 and length is 2
14 starts at 19 and length is 2

```

下面是使用命名分组的程序（结果相同，略去）。

```

Regex regex = new Regex(@"(?<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})");
Match match = regex.Match("2010-12-20 2011-02-14");
String[] groupNames = regex.GetGroupNames();
while (match.Success) {
    for (int i = 0, n = groupNames.Length; i < n; i++) {
        String groupName = groupNames[i];
        Console.WriteLine(match.Groups[groupName].Value + " starts at " +
            match.Groups[groupName].Index + " and length is " +
            match.Groups[groupName].Length);
    }
    match = match.NextMatch();
}

```

注意：如果在 VB.NET 中访问 Match 对象的分组，不能使用 Groups[num] 或 Groups[name] 之类的方法，而应该使用 Groups(num) 或 Groups(name)。

## 10.4 常用操作示例

### 10.4.1 验证

#### 简单验证

```

//注意要加上\A 和\z
String regex = @"\A\d{4}-\d{2}-\d{2}\z";
Regex.IsMatch("2010-12-20", regex); // => True

```

#### 用 Regex 对象验证，方便反复使用

```

//同样要加上\A 和\z
Regex regex = new Regex(@"\A\d{4}-\d{2}-\d{2}\z");
String s = "2010-12-20";
regex.IsMatch(s); // => True

```

### 10.4.2 提取

#### 使用数字编号分组进行数据提取

```

Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
Match match = regex.Match("2010-12-20 2011-02-14");
int[] groupNums = regex.GetGroupNumbers();

```

```

while (match.Success) {
    Console.Write("date: " + match.Groups[0].Value);
    Console.Write(", year: " + match.Groups[1].Value);
    Console.Write(", month: " + match.Groups[2].Value);
    Console.Write(", day: " + match.Groups[3].Value);
    Console.Write("\n");
    match = match.NextMatch();
}
date: 2010-12-20, year: 2010, month: 12, day: 20
date: 2011-02-14, year: 2011, month: 02, day: 14

```

#### 使用命名分组进行数据提取

```

Regex regex = new Regex(@"(<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})");
Match match = regex.Match("2010-12-20 2011-02-14");
int[] groupNums = regex.GetGroupNumbers();
while (match.Success) {
    Console.Write("date: " + match.Groups["0"].Value);
    Console.Write(", year: " + match.Groups["year"].Value);
    Console.Write(", month: " + match.Groups["month"].Value);
    Console.Write(", day: " + match.Groups["day"].Value);
    Console.Write("\n");
    match = match.NextMatch();
}
date: 2010-12-20, year: 2010, month: 12, day: 20
date: 2011-02-14, year: 2011, month: 02, day: 14

```

### 10.4.3 替换

#### 简单替换

```

String regex = @"(\d{4})-(\d{2})-(\d{2})";
String replacement = "$2/$3/$1";
String text = "2010-12-20 2011-02-14";
Console.WriteLine(Regex.Replace(text, regex, replacement));
12/20/2010 02/14/2011

```

#### 用 Regex 对象替换，方便反复使用

```

Regex regex = new Regex(@"(\d{4})-(\d{2})-(\d{2})");
String replacement = "$2/$3/$1";
String text = "2010-12-20 2011-02-14";
Console.WriteLine(regex.Replace(text, replacement));
12/20/2010 02/14/2011

```

#### 使用命名分组

```

String regex = @"(<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})";
String replacement = @"${month}/${day}/${year}";

```

```
String text = "2010-12-20 2011-02-14";  
Console.WriteLine(Regex.Replace(text, regex, replacement));  
12/20/2010 02/14/2011
```

在 replacement 中使用 \$

```
Regex.Replace("the price is 12.99", @"\d+\.\d{0,2}", "$$$0");  
the price is $12.99
```

## 10.4.4 切分

简单切分

```
String regex = @"-";  
foreach(String s in Regex.Split("2010-12-20", regex)) {  
    Console.Write(s + " ");  
}  
2010 12 20
```

用 Regex 对象切分，方便重复使用

```
Regex regex = new Regex("-");  
foreach(String s in regex.Split("2010-12-20")) {  
    Console.Write(s + " ");  
}  
2010 12 20
```

# 第 11 章 Java

Java 对正则表达式的支持比较完备，在 JDK 1.4 之后，Java 内置了 `java.util.regex` 包，提供了几乎所有常用的功能；Java 5（Java 1.5）和 Java 6（Java 1.6）在正则表达式的使用上保持了一贯，主要改进是修正了少数错误。考虑到如今 Java 5 是普遍使用的版本，本章的讲解以 Java 5 为准<sup>1</sup>。本章介绍的内容也完全适用于 Android 开发，Android 程序中正则表达式的使用方法本文所述的相同。

值得一提的是，Java 语言中的字符都是 Unicode 编码的，Java 的正则表达式对 Unicode 的兼容也非常好。

## 11.1 预备知识

Java 语言中的正则表达式相关类都存在于 `java.util.regex` 包中，如果要使用正则表达式，必须首先导入其中的类，最简单的办法是在程序文件头部写入这条语句（具体的类在后文详述）。

```
import java.util.regex.*;
```

为讲解方便，先介绍 Java 中的 `String.matches(String regex)` 方法，它返回一个布尔值，表示整个 `String` 能不能由 `regex` 匹配。

```
"1".matches("\\d");    // => True
"1a".matches("\\d");   // => False
```

如果要检查正则表达式能否在字符串中找到匹配（也就是匹配到某个子串），可以使用下面的方法（为排版方便，这里用到了链式编程，关于 `Pattern` 的具体细节在后面详述）。

```
Pattern.compile("\\d").matcher("1a").find(); // => True
```

Java 语言中的正则表达式都是以字符串的形式给出的，所以字符串转义和正则表达式转义都必须考虑。如果正则表达式中出现了反斜线 `\`，在通过字符串给出时，都必须进行字符串转义，写为 `\\`。比如正则表达式中的 `\b`，在字符串中应该写作 `\\b`，`\Q...\E` 则应当写作 `\\Q...\\E`。稍微麻烦一点的是，如果要在正则表达式中使用反斜线 `\`，它本身应当转义为 `\\`，而在字符串中，两

---

<sup>1</sup> 更详细的文档可见于：

<http://download.oracle.com/javase/1.5.0/docs/api/java/util/regex/Pattern.html>

<http://download.oracle.com/javase/1.5.0/docs/api/java/util/regex/Matcher.html>

个反斜线应当分别转义，所以字符串中应该写作\\\\"。

```
//这个字符串只包含反斜线字符\  
"\\\".matches("\\\\"); // => True
```

## 11.2 正则功能详解

### 11.2.1 列表

表 11-1 中列出了 Java 中的正则功能。

表 11-1 Java 中的正则功能列表

功能	记法	说明
字符组 ④1	<code>[...]</code> <code>[^...]</code>	完全支持 Unicode
POSIX 字符组 ④15	<code>[:digit:]</code>	只匹配 ASCII 字符，写成 <code>\P{...}</code>
Unicode 属性 ④121	<code>\P{...}</code>	完全支持 Unicode
字符组简记法 ④115	<code>\d</code> <code>\D</code> <code>\w</code> <code>\W</code> <code>\s</code> <code>\S</code>	采用 ASCII 匹配规则
行起始位置 ④61	<code>^</code> <code>\A</code>	支持
行结束位置 ④61	<code>\$</code> <code>\z</code> <code>\Z</code>	支持
单词边界 ④117	<code>\b</code> <code>\B</code>	采用 Unicode 规则
顺序环视 ④68	<code>(?=...)</code> <code>(?!...)</code>	支持
逆序环视 ④68	<code>(?&lt;=...)</code> <code>(?&lt;!...)</code>	不支持无法确定长度范围的表达式
匹配模式 ④81	<code>i</code> <code>m</code> <code>s</code> <code>x</code>	支持
模式作用范围 ④89	<code>(?modifier)</code> <code>(?-modifier)</code>	支持
纯文本模式 ④99	<code>\Q...\E</code>	支持(存在小问题,到 Java 6 才修正)
捕获分组及引用④44	<code>(...)</code> <code>\num</code> <code>\$num</code>	支持
命名分组 ④53	<code>(?&lt;name&gt;...)</code> <code>\k{name}</code> <code>\${name}</code>	不支持 (Java 7 支持)
非捕获分组 ④54	<code>(?:...)</code>	支持
多选结构 ④39	<code>(... ...)</code>	支持
匹配优先量词 ④18	<code>?</code> <code>*</code> <code>+</code>	支持
忽略优先量词 ④26	<code>??</code> <code>?*</code> <code>?+</code> <code>{n,m}?</code>	支持

### 11.2.2 字符组

在 Java 语言中所有的字符都采用 UTF-16 编码，所以在字符组中可以直接使用中文，因而不会出现多字节字符错误匹配的问题，而且点号 `.` 可以匹配中文字符。

```
"你".matches("."); // => True
"罚".matches("[收发]"); // => False
```

值得注意的是，Java 语言中的字符组可以进行集合运算，通常主要用于字符组的“减”运算——如果需要匹配英文中所有小写的辅音字母，也就是从 26 个小写字母中“减去”5 个元音字母，可以写作 `[[a-z]&&[^aeiou]]`。它的意思是，“从 a 到 z 的所有字符（也就是小写英文字符）”与“除 a、e、i、o、u 之外的任何字符”取交集，也就是就是所有的小写辅音字母。

```
"a".matches("[a-z]&&[^aeiou]"); // => False
"b".matches("[a-z]&&[^aeiou]"); // => True
```

也可以用这个功能实现匹配“英文大小写字母和数字（不包括下划线）”。

```
"a".matches("[\\w&&[^_] ]"); // => True
"0".matches("[\\w&&[^_] ]"); // => True
"_" .matches("[\\w&&[^_] ]"); // => False
```

相对于 `&&` 的“交”运算，Java 也支持在字符组内用 `|` 进行“并”运算，不过因为“并集”的意思就是“给字符组内部添加字符”，所以并不需要设定特殊的运算符，比如 `[[0-4]||[6-9]]` 就等价于 `[[0-4][6-9]]`，也等价于 `[0-46-9]`。虽然这几种写法都没错，但最后那种写法显然更简单，所以我推荐最后的写法。

```
"3".matches("[[0-4][6-9]]"); // => True
"5".matches("[[0-4][6-9]]"); // => False
```

因为“并”是默认的运算，如果字符组内部只出现一个字符组，可以理解为“单个字符组与空集取并集”，换句话说，`[[0-4]]` 等价于 `[0-4]`。看到 `[[0-4]]`，而不是包含七个字符（0、1、2、3、4、[、]）的字符组。

```
"3".matches("[[0-4]]"); // => True
"[".matches("[[0-4]]"); // => False
```

Java 中可以用 `\uhex` 指定 Unicode 码值，其中 `hex` 为 4 位十六进制数，所以在 Unicode 编码环境下可以用字符组 `[\u4e00-\u9fff]` 匹配所有的中文字符。

```
"我".matches("[\u4e00-\u9fff]"); // => True
```

Java 中也可以使用 POSIX 字符组 `[:name:]`，不过必须写为 `\p{name}`，而且第一个字母必须大写，比如 `[:digit:]` 就应该写成 `\p{Digit}`，其他 POSIX 字符组都是如此，只有 `[:xdigit:]` 必须写成 `\p{XDigit}`。注意，所有的 POSIX 字符组都只能匹配 ASCII 字符。

```
"0".matches("\\p{Digit}"); // => True
"f".matches("\\p{XDigit}"); // => True
```

### 11.2.3 Unicode 属性

Java 对 Unicode 字符组的支持比较好，它支持 Unicode Property。这样，用 `\p{N}` 就可以匹配所有的数字字符，包括中文的全角数字字符，比如 0、1、2 等；`\p{P}` 可以匹配各种标点符号，包括中文的冒号：等。

```
"1".matches("\\p{N}"); // => True
", ".matches("\\p{P}"); // => True
```

Java 也支持 Unicode Block，其记法是以 In 为前缀的，比如匹配中文字符的 Unicode Block 就应当写为 `InCJK_Compatibility_Ideographs`，这与 .NET 不一样，使用时应注意<sup>1</sup>。

```
"我".matches("\\p{InCJK_Compatibility_Ideographs}"); // => True
```

关于 Unicode 属性的细节，请参考第 123 页。

### 11.2.4 字符组简记法

在 Java 的正则表达式中，`\d`、`\D`、`\w`、`\W`、`\s`、`\S` 都使用 ASCII 匹配规则。也就是说，`\d` 等价于 `[0-9]`，`\w` 等价于 `[0-9a-zA-Z_]`，`\s` 无法匹配 ASCII 编码之外的空白字符。

```
//半角字符
"1".matches("\\d"); // => True
"a".matches("\\w"); // => True
" ".matches("\\s"); // => True
//全角字符及中文字符
"1".matches("\\d"); // => False
"我".matches("\\w"); // => False
" ".matches("\\s"); // => False
```

### 11.2.5 单词边界

Java 正则表达式中的单词边界有点特殊：一般来说，`\b` 的匹配是受到 `\w` 影响的，但 Java 却不是如此（☞ 117）。尽管 `\w` 使用 ASCII 规则，`\b` 定义的“单词字符”是遵循 Unicode 规则的。

这一点官方文档并没有详细解释，不过我们可以这样粗略记忆：`\b` 所识别的单词字符，可以是各语言中的单词字符（包括中文字符），但不包括空白和标点（包括中文的空白和标点）；而且，单词字符必须出现。

```
//英文字符-半角标点
Pattern.compile("a\\b").matcher("a,").find(); // => True
```

---

<sup>1</sup> Unicode Block 的细节可参考 <http://download.oracle.com/javase/6/docs/api/java/lang/Character.UnicodeBlock.html>。

```
//英文字符-全角标点
Pattern.compile("a\\b").matcher("a, ").find(); // => True
//英文字符-结束
Pattern.compile("a\\b").matcher("a ").find(); // => True
//中文字符-半角标点
Pattern.compile("我\\b").matcher("我, ").find(); // => True
//中文字符-全角标点
Pattern.compile("我\\b").matcher("我, ").find(); // => True
//中文字符-空字符串
Pattern.compile("我\\b").matcher("我").find(); // => True
//全角标点-结束（未出现单词字符）
Pattern.compile(", \\b").matcher(", ").find(); // => False
//中文字符-英文字符（未出现单词字符）
Pattern.compile("我\\b").matcher("我 a").find(); // => False
```

### 11.2.6 行起始/结束位置

^匹配的是“行开始的位置”，在默认情况下它只能匹配“整个字符串的开始位置”，如果指定使用多行模式（Multiline Mode），^可以匹配字符串内部的文本行的开始位置；而^A无论在什么情况下都匹配“整个字符串的开始位置”。

```
Pattern.compile("^1").matcher("1\n2\n").find(); // => True
Pattern.compile("^2").matcher("1\n2\n").find(); // => False
Pattern.compile("(?m)^2").matcher("1\n2\n").find(); // => True

Pattern.compile("^\\A1").matcher("1\n2\n").find(); // => True
Pattern.compile("^\\A2").matcher("1\n2\n").find(); // => False
Pattern.compile("(?m)^\\A2").matcher("1\n2\n").find(); // => False
```

在默认情况下，Java 中的\$、\\Z、\\z匹配的是整个字符串的结束位置（如果结束位置有换行符，则\$和\\z匹配这个换行符之前的位置），多行模式只会影响到\$的匹配，这时候它可以匹配文本内部行的结束位置，所以进行准确的验证时应当使用\\z。

```
Pattern.compile("1$").matcher("1\n").find(); // => True
Pattern.compile("1$").matcher("1\n2").find(); // => False
Pattern.compile("(?m)1$").matcher("1\n2").find(); // => True

Pattern.compile("1\\Z").matcher("1\n").find(); // => True
Pattern.compile("1\\Z").matcher("1\n2").find(); // => False
Pattern.compile("(?m)1\\Z").matcher("1\n").find(); // => True

Pattern.compile("1\\z").matcher("1").find(); // => True
Pattern.compile("1\\z").matcher("1\n").find(); // => False
Pattern.compile("(?m)1\\z").matcher("1\n").find(); // => False
```



11.2.7 环视

Java 中的肯定顺序环视 `(?=...)` 和否定顺序环视 `(?!...)` 内可以使用任意形式的子表达式。

```
"cab".matches("c(?=ab+)");           // => True
"ccd".matches("c(?(=ab+|cd))");       // => True
"cabbbb".matches("c(?(=ab+|cd))");    // => True
```

逆序环视中使用的子表达式则**必须有长度上限**。也就是说，`(?<=ab)`、`(?<=ab?)`、`(?<=ab{4,12})` 都是合法的，`(?<=ab+)` 则会报错。这个问题值得注意，因为它在代码编译阶段不会抛出异常，而在运行时会出现。

```
"ab".matches("ab(?<=ab)");           // => True
"cd".matches("cd+(?<=(abcd|cd))");    // => True
"cd".matches("cd(?<=(ab+|cd))");    //编译时没有问题，运行却会报错
```

11.2.8 匹配模式

表 11-2 列出了 Java 中的正则表达式支持的匹配模式。

表 11-2 Java 中的正则表达式支持的匹配模式

常量	修饰符	说明
CASE_INSENSITIVE	i	不区分 ASCII 字符的大小写
COMMENTS	x	容许正则表达式中出现注释，表达式中的空白字符，以及#开始到行末的文本，都视为注释
MULTILINE	m	容许 <code>^</code> 和 <code>\$</code> 不仅仅匹配字符串的起始和结束位置，还可以匹配字符串内部文本行的起始和结束位置
DOTALL	s	容许点号 <code>.</code> 匹配任何字符，包括换行符
UNICODE_CASE	u	可以识别 Unicode 字符的不同形态，“不区分大小写”的范围不限于 ASCII 字符，但严重影响性能
UNIX_LINES	d	限定 <code>.</code> 、 <code>^</code> 、 <code>\$</code> 能识别的行终结符只有换行符 <code>\n</code> ，忽略 <code>\r\n</code> 和 Unicode 行终结符等其他字符
CANON_EQ	无	匹配时采取 Unicode “等价”规则，可以识别意义相等的复合字符（单个字符加上调号）与单个字符；但此选项会严重影响性能

匹配模式可以在表达式中以 `(?modifier)` 指定，也可以在生成 `Pattern` 对象时，以预定义常量指定。

```
"A".matches("a");           // => False
"A".matches("(?i)a");       // => True
Pattern.compile("a", Pattern.CASE_INSENSITIVE).matcher("A").find(); // => True
```

Java 也支持用 `(?-modifier)` 停用某个匹配模式。

```
"aBb".matches("a(?i)b(?-i)b");    // => True
"aBB".matches("a(?i)b(?-i)b");    // => False
```

### 11.2.9 纯文本模式

在 `\Q` 和 `\E` 之内，所有的元字符和特殊结构都失去特殊意义，只能匹配它们对应的字符本身。不过在字符串中，`\Q` 和 `\E` 中的反斜线也需要经过字符串转义，应该写成 `\\Q` 和 `\\E`<sup>1</sup>。

```
"A".matches("\\Q.\\E");            // => False
"." .matches("\\Q.\\E");            // => True
"[a-f]".matches("\\Q[a-f]\\E");    // => True
"b".matches("\\Q[a-f]\\E");        // => False
```

### 11.2.10 捕获分组的引用

在 Java 中，如果要在正则表达式内部引用捕获分组，应当使用 `\num` 记法，其中 *num* 为对应捕获分组的编号。

```
"ab".matches("([a-z])\\1");        // => False
"aa".matches("([a-z])\\1");        // => True
```

要在替换时引用捕获分组，应当使用 `$num` 记法。

```
"2010-12-20".replaceAll("(\\d{4})-(\\d{2})-(\\d{2})", "$2/$3/$1");
12/20/2010
```

要在 *replacement* 字符串中表示 `$` 字符，必须使用 `\\$` 转义。

```
"the price is 12.99".replaceAll("\\d+\\.\\d{0,2}", "\\$0");
the price is $12.99
```

## 11.3 正则 API 简介

Java 中正则表达式相关的类都在 `java.util.regex` 之内，一般来说，主要用到的是这两个类：`java.util.regex.Pattern`（以下简称 *Pattern*）和 `java.util.regex.Matcher`（以下简称 *Matcher*）。*Pattern* 对应正则表达式，一个 *Pattern* 与一个 *String* 对象关联，生成一个 *Matcher*，它对应 *Pattern* 在 *String* 中的一次匹配；调用 *Matcher* 对象的 `find()` 方法，*Matcher* 对象就会更新为下一次匹配的匹配信息。*Pattern*、*String*、*Matcher* 对象之间的联系如图 11-1 所示。

---

<sup>1</sup> 但是此功能在 Java 5 中有 bug，如果在字符组内部使用，可能会出现奇怪的问题，到了 Java 6 才可以彻底解决。

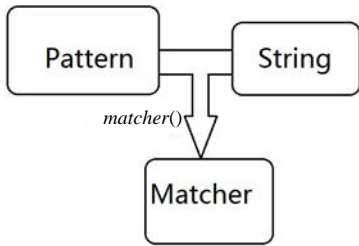


图 11-1 Java 中与正则表达式相关的类

## 例 11-1 Java 中正则表达式相关类的实际使用

```

Pattern pattern = Pattern.compile("\\d{4}-\\d{2}-\\d{2}");
String string = "2010-12-20 2011-02-14";
Matcher matcher = pattern.matcher(string);
while(matcher.find()) {
    System.out.println(matcher.group(0));
}

```

2010-12-20

2011-02-14

## 11.3.1 Pattern

Pattern 是 Java 语言中的正则表达式对象（字符串与正则表达式的差别详见第 8 章）。要使用正则表达式，首先必须从字符串“编译”出 Pattern 对象，这需要用到 `Pattern.compile (String regex)` 方法。

```
Pattern pattern = Pattern.compile("a.b+");
```

如果要指定匹配模式，可以在表达式中使用 `(?modifier)` 修饰符指定，也可以使用预定义常量。下面的两个 Pattern 对象的生成方法不同，结果却是等价的。

```

Pattern pattern = Pattern.compile("(?i)a.b+");
Pattern pattern = Pattern.compile("a.b+", Pattern.CASE_INSENSITIVE);

```

如果要同时指定多种模式，可以连写模式修饰符，也可以直接用 `|` 运算符将预定义常量连接起来，以下两个 Pattern 对象也是等价的。

```

Pattern pattern = Pattern.compile("(?is)a.b+");
Pattern pattern = Pattern.compile("a.b+", Pattern.CASE_INSENSITIVE | Pattern.DOTALL);

```

下面介绍 Pattern 的主要成员方法。

### 11.3.1.1 static boolean matches(String regex, CharSequence input)

这个方法可以检验字符串 *input* 能否由正则表达式 *regex* 匹配，因为是静态方法，所以不需要编译生成各个对象，方便随手使用。要注意的是，它检验的是“整个字符串能否由表达式匹配”，而不是“表达式能否在字符串中找到匹配”。你可以认为 *regex* 的首尾自动加上了匹配字符串起始和结束位置的锚点 `\A` 和 `\Z`。

```
Pattern.matches("\\d{6}", "a123456"); // => False
Pattern.matches("\\d{6}", "123456"); // => True
```

### 11.3.1.2 String[] split(CharSequence text)

通常，`Pattern` 对象需要配合下面将要介绍的 `Matcher` 一起完成正则操作。如果只用正则表达式来切分字符串，只用 `Pattern` 的这个方法也可以。

这个方法接收的参数类型是 `CharSequence` 它可能有点陌生，其实它是 `String` 的父类，其他子类还有 `CharBuffer`、`StringBuffer`、`StringBuilder`，因而可以应对常见的各种表示“字符串”的类。下面的代码仅以 `String` 为例。

```
String s = "2010-12-20-";
Pattern pattern = Pattern.compile("\\s+");
for(String part : pattern.split(s)) {
    System.out.println(part);
}
2010
12
20
```

请注意，虽然正则表达式能在字符串中匹配 3 次，但切分得到的结果数是 3 个而不是 4 个，因为末尾空白字符之后的空字符串会被自动忽略。如果这个空字符串出现在其他位置（比如字符串开头，或字符串中），则会被保留。

### 11.3.1.3 String[] split(CharSequence text, int limit)

这个方法与上面的方法很相似，只是多了一个参数 *limit*，它用来限定返回的 `String` 数组的最大长度。也就是说，它规定了字符串至多只能“切”*limit*-1 次。如果不需要对字符串比较大，进行尽可能多的切分，使用这个方法。

```
String s = " 2010-12-20 ";
Pattern pattern = Pattern.compile("\\s+");
for(String part : pattern.split(s, 2)) {
    System.out.println(part);
}
2010
12-20
```

既然 *limit* 是一个 `int` 类型，那么它自然可以设定为各种值，表 11-3 总结了 *limit* 在各个取值

区间对结果的影响（未指定 *limit* 时，最终返回包含 *n* 个元素的数组，实际能切分的次数是 *n-1*）。

表 11-3 limit 取值对结果的影响

取值	结果
limit<0	等于未设定 limit 时，保留末尾的空字符串
limit=0	等于未设定 limit 时，切分 <i>n-1</i> 次，忽略末尾的空字符串
0 <limit<n	返回数组包含 limit 个元素，切分 limit-1 次，最后一个元素是第 limit-1 次切分后，右侧剩下的所有文本
limit>=n	等于未指定 limit 时

11.3.1.4 static String quote(String text)

这个方法用来取消字符串 *text* 中所有转义字符的特殊含义，实质就是在字符串首尾添加 `\Q` 和 `\E`。通常，如果需要把某个字符串作为没有任何特殊意义的正则表达式（比如从外界读入的字符串，用在某个复杂的正则表达式中），就可以使用这个方法。

```
"acb".matches("a*.b"); // => True
"a*.b".matches("a*.b"); // => False
"a*.b".matches("a*.b"); // => False
"a*.b".matches(Pattern.quote("a*.b")); // => True
```

11.3.2 Matcher

Matcher 可以理解为“某次具体匹配的结果对象”：把编译好的 Pattern 对象“应用”到某个 String 对象上，就获得了作为“本次匹配结果”的 Matcher 对象。之后，就可以通过它获得关于匹配的信息。

```
Pattern pattern = Pattern.compile("\\d{4}-\\d{2}-\\d{2}");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
while(matcher.find()) {
    System.out.println(matcher.group());
}
2010-12-20
2011-02-14
```

对编译好的 Pattern 对象调用 `matcher(String text)` 方法，传入要匹配的字符串 *text*，就得到了 Matcher 对象，每调用一次 `find()` 方法，如果返回 `true`，就表示“找到一个匹配”，此时可以通过下面的若干方法获得关于本次匹配的信息。

11.3.2.1 String group(int n)

返回当前匹配中第 *n* 对捕获括号捕获的文本，如果 *n* 为 0，则取匹配的全部内容；如果 *n* 小于 0 或者大于最大分组编号数，则报错。

### 11.3.2.2 String group()

返回当前匹配的全部文本，相当于 `group(0)`。

### 11.3.2.3 int groupCount()

返回此 `Matcher` 对应 `Pattern` 对象中包含的捕获分组数目，编号为 0 的默认分组不计在内。

### 11.3.2.4 int start(n)

返回当前匹配中第 `n` 对捕获括号匹配的文本在原字符串中的起始位置。

### 11.3.2.5 int start()

返回当前匹配的文本在原字符串中的起始位置，相当于 `start(0)`。

### 11.3.2.6 int end(n)

返回当前匹配中第 `n` 对捕获括号匹配的文本在原字符串中的结束位置。

### 11.3.2.7 int end()

返回当前匹配的文本在原字符串中的结束位置，相当于 `end(0)`。

下面的例子集中示范这几个方法。

```
Pattern pattern = Pattern.compile("(\\d{4})-(\\d{2})-(\\d{2})");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
while(matcher.find()) {
    System.out.println(matcher.group() + " starts at " + matcher.start() + " and end
at " + matcher.end());
    for (int i = 0, n = matcher.groupCount(); i <= n; i++) {
        //注意这里是 i<=n, 因为编号为 n 的分组一定存在
        System.out.println(matcher.group(i) + " starts at " + matcher.start(i) + " and
ends at " + matcher.end(i));
    }
}
```

```
2010-12-20 starts at 0 and ends at 10
2010-12-20 starts at 0 and ends at 10
2010 starts at 0 and ends at 4
12 starts at 5 and ends at 7
20 starts at 8 and ends at 10
2011-02-14 starts at 11 and ends at 21
2011 starts at 11 and ends at 15
02 starts at 16 and ends at 18
14 starts at 19 and ends at 21
```

### 11.3.2.8 String replaceAll(String replacement)

如果进行正则表达式替换，一般用到的是 `Matcher` 的 `replaceAll()` 方法，它会将原有文本中正则表达式能匹配的所有文本替换为 *replacement* 字符串。

```
Pattern pattern = Pattern.compile("\\d{4}-\\d{2}-\\d{2}");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
System.out.println(matcher.replaceAll("Date"));
Date Date
```

如果要在 *replacement* 字符串中引用之前某个捕获分组匹配的文本，则应当使用 *\$num*，其中 *num* 为捕获分组的编号。

```
Pattern pattern = Pattern.compile("(\\d{4})-(\\d{2})-(\\d{2})");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
System.out.println(matcher.replaceAll("$2/$3/$1"));
12/20/2010 02/14/2011
```

如果要在 *replacement* 中表示 `$` 符号，则应当将它转义为 `\$`。

```
Pattern pattern = Pattern.compile("(\\d+\\.\\d+)");
Matcher matcher = pattern.matcher("the price is 12.99");
System.out.println(matcher.replaceAll("\\$0"));
the price is $12.99
```

## 11.3.3 String

许多时候只需要临时使用某个正则表达式，而不需要重复使用，这时候每次都生成 `Pattern` 对象和 `Matcher` 对象再操作显得很烦琐。所以，Java 的 `String` 类提供了正则表达式操作的静态成员方法，只需要 `String` 对象就可以执行正则表达式操作。下面讲解 `String` 类常见的正则操作。

### 11.3.3.1 boolean matches(String regex)

这个方法判断当前的 `String` 对象能否由正则表达式 *regex* 匹配。请注意，这里的“匹配”指的并不是 *regex* 能否在 `String` 内找到匹配，而是指 *regex* 匹配整个 `String` 对象，因此非常适合用来做数据校验。

```
"123456".matches("\\d{6}"); // => True
"a123456".matches("\\d{6}"); // => False
```

### 11.3.3.2 String replaceFirst(String regex, String replacement)

这个方法用来替换正则表达式 *regex* 在字符串中第一次能匹配的文本，可以在 *replacement* 字符串中用 *\$num* 引用 *regex* 中对应捕获分组匹配的文本。

```
"2010-12-20 2011-02-14".replaceFirst("\\d{4}-\\d{2}-\\d{2}", "Date");
Date 2010-02-14
```

```
"2010-12-20 2011-02-14".replaceFirst("(\\d{4})-(\\d{2})-(\\d{2})", "$2/$3/$1");
12/20/2010 2011-02-14
```

### 11.3.3.3 String replaceAll(String regex, String replacement)

这个方法用来进行所有的替换，它的结果等同于 `Matcher` 类的 `replaceAll()` 方法，`replacement` 字符串中也可以用 `$num` 的表示法引用 `regex` 中对应捕获分组匹配的文本。

```
"2010-12-20 2011-02-14".replaceAll("\\d{4}-\\d{2}-\\d{2}", "Date");
Date Date

"2010-12-20 2011-02-14".replaceAll("(\\d{4})-(\\d{2})-(\\d{2})", "$2/$3/$1");
12/20/2010 02/14/2011
```

值得注意的是，在 Java 5.0 之前的版本中，`String` 类并没有提供简单的“字符串替换”方法，只有 `replaceAll()`。如果你使用的 Java 版本低于 5.0，又需要进行简单的字符串替换，可以使用 `replaceAll()`，不过应该在要被替换的字符串首尾加上 `\Q` 和 `\E`。

### 11.3.3.4 String[] split(String regex)

这个方法等价于 `Pattern` 中对应的 `split()` 方法，此处不再赘述。

### 11.3.3.5 String[] split(String regex, int limit)

这个方法等价于 `Pattern` 中对应的 `split()` 方法，此处不再赘述。

`String` 类的这些方法，其实都是包装了 `Pattern` 对象和 `Matcher` 对象，只是操作更方便，并不能提高运行的效率。

## 11.4 常用操作示例

### 11.4.1 验证

简单验证

```
"2010-12-20".matches("\\d{4}-\\d{2}-\\d{2}"); // => True
Pattern.matches("\\d{4}-\\d{2}-\\d{2}", "2010-12-20"); // => True
```

使用 `Pattern` 对象，方便反复验证

```
//注意要加上\\A和\\z
Pattern pattern = Pattern.compile("\\A\\d{4}-\\d{2}-\\d{2}\\z");
String s = "2010-12-20";
pattern.matcher(s).find(); // => True
```



### 11.4.2 提取

```
Pattern pattern = Pattern.compile("\\d{4}-\\d{2}-\\d{2}");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
while(matcher.find()) {
    System.out.println(matcher.group());
}
2010-12-20
2011-02-14
```

```
Pattern pattern = Pattern.compile("(\\d{4})-(\\d{2})-(\\d{2})");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
while(matcher.find()) {
    System.out.print("date: " + matcher.group(0));
    System.out.print(", year: " + matcher.group(1));
    System.out.print(", month: " + matcher.group(2));
    System.out.print(", day: " + matcher.group(3));
    System.out.print("\n");
}
date: 2010-12-20, year: 2010, month: 12, day: 20
date: 2011-02-14, year: 2011, month: 02, day: 14
```

### 11.4.3 替换

#### 简单替换

```
String regex = "(\\d{4})-(\\d{2})-(\\d{2})";
String replacement = "$2/$3/$1";
"2010-12-20 2011-02-14".replaceAll(regex, replacement);
12/20/2010 02/14/2011
```

#### 使用 Pattern 对象，方便反复使用

```
Pattern pattern = Pattern.compile("(\\d{4})-(\\d{2})-(\\d{2})");
Matcher matcher = pattern.matcher("2010-12-20 2011-02-14");
System.out.println(matcher.replaceAll("$2/$3/$1"));
12/20/2010 02/14/2011
```

#### 在 replacement 中使用 \$

```
String regex = "\\d+\\.\\d{0,2}";
String replacement = "\\$0";
System.out.println("price is 12.99".replaceAll(regex, replacement));
price is $12.99
```

### 11.4.4 切分

#### 简单切分

```
String regex = "\\s+";
for(String s : "one\ttwo\n three".split(regex)) {
    System.out.print(s + " ");
}
one two three
```

#### 使用 Pattern 对象，方便反复使用

```
Pattern pattern = Pattern.compile("\\s+");
for(String s : pattern.split("one\ttwo\n three")) {
    System.out.println(s);
}
one two three
```

## 第12章 JavaScript

JavaScript 正则表达式的情况比较复杂，因为同样的 JavaScript 代码在不同浏览器上的表现可能不同，幸运的是，本章介绍的功能几乎都是“通用”的，特别情况会专门标注出来。

JavaScript 1.2（发布于1997年）以后的版本都内建了正则表达式支持，从 JavaScript 1.5（发布于2000年11月）开始，对正则表达式的支持更加完备了。JavaScript 1.5 的规范对应于 ECMA-262 Edition 3 的标准，目前流行的所有主流浏览器（包括 IE、Firefox、Chrome、Safari 等）都支持此标准。所以通常来说，JavaScript 的正则表达式在各种浏览器里的表现都是相同的，但是也有些例外。如果同样的表达式在不同的浏览器中表现不同，本章会注明。

Adobe Flash 所使用的 ActionScript 也遵循 ECMA-262 标准，所以，JavaScript 的正则表达式可以很容易地使用在 ActionScript 程序中，本章末尾有专门介绍 ActionScript 的内容。

### 12.1 预备知识

JavaScript 的正则表达式对象是 `RegExp`<sup>1</sup>。最常见的生成方法是在首尾两个斜线分隔符/之间，用正则文字给出表达式，比如 `/\d+/`。

为了讲解方便，现在先介绍 `RegExp` 的 `test()` 方法：`RegExp.test(string)` 返回一个布尔值，表示表达式能否在 `string` 中找到匹配。注意它并不要求整个 `string` 由表达式匹配，如果要检验整个 `string` 能否由表达式匹配，可以在表达式两端添加 `^` 和 `$`，因为 JavaScript 不支持 `\A`、`\Z`、`\z`。

```
/\d/.test("1");           // => True
/\d/.test("a1");          // => True
/^\.test("a1");           // => False
```

也可以显式生成 `RegExp`，有两种办法：一种是使用字符串文字，这时需要考虑字符串转义和正则表达式转义；另一种是使用正则文字，直接在表达式两端写上反斜线/作为分隔符即可。这时候不需要考虑字符串的转义，唯一需要特殊考虑的字符就是反斜线/，如果它在正则表达式内部出现，必须转义为 `\/`。相比起来，正则文字的方式更直观，所以如果要显式生成 `RegExp` 对象，推荐使用后一种方法。

---

<sup>1</sup> 如果你用过 Ruby 中的正则表达式，可能会发现 Ruby 和 JavaScript 中的正则表达式对象名字很像，唯一的区别在于 Ruby 的 `Regexp` 中 `e` 是小写，而 JavaScript 的 `RegExp` 中 `E` 是大写。

```
(new RegExp("\\d")).test("1");    // => True
(new RegExp(/\d/)).test("1");    // => True
(new RegExp(/\d/)).test("/");    // => True
```

## 12.2 正则功能详解

### 12.2.1 列表

表 12-1 中列出了 JavaScript 中的正则功能。

表 12-1 JavaScript 中的正则功能

功能	记法	说明
字符组 ④1	<code>[...] [^...]</code>	可以使用 Unicode 字符
POSIX 字符组 ④15	<code>[:digit:]</code>	不支持
Unicode 属性 ④121	<code>\p{...}</code>	不支持
字符组简记法 ④115	<code>\d \D \w \W \s \S</code>	采用 ASCII 匹配规则
行起始位置 ④61	<code>^</code>	只支持 <code>^</code>
行结束位置 ④61	<code>\$</code>	只支持 <code>\$</code>
单词边界 ④117	<code>\b \B</code>	采用 ASCII 匹配规则
顺序环视 ④68	<code>(?=...) (?!...)</code>	支持（限 JavaScript 1.5 以后）
逆序环视 ④68	<code>(?&lt;=...) (?&lt;!...)</code>	不支持
匹配模式 ④81	<code>i m g</code>	常用模式中只支持 <code>i</code> 和 <code>m</code> , <code>g</code> 模式在下文详细讲解
模式作用范围 ④89	<code>(?modifier) (?~modifier)</code>	不支持
纯文本模式 ④99	<code>\Q...\E</code>	不支持
捕获分组及引用 ④44	<code>(...) \num \$num</code>	支持
命名分组 ④53	<code>(?&lt;name&gt;...) \k{name} \${name}</code>	不支持
非捕获分组 ④54	<code>(?:...)</code>	支持
多选结构 ④39	<code>(... ...)</code>	支持
匹配优先量词 ④18	<code>? * +</code>	支持
忽略优先量词 ④26	<code>?? ?* ?+ {n,m}?</code>	支持

### 12.2.2 字符组

在 JavaScript 的正则表达式中，中文字符问题比较复杂，如果 JavaScript 程序使用的是非 Unicode 编码（比如常见的 GBK 编码），在某些浏览器中，下面的程序可能返回 `true`（具体取决

于所使用的浏览器)。

```
[/收发]/.test("罚");    // => True
```

解决办法是使用 Unicode 编码环境,这样正则表达式不会将多字节字符拆成“字节”来对待,可以避免错误匹配。其中具体的细节,请参考第 111 页的讲解。

```
[/收发]/.test("罚");    // => False
```

JavaScript 中可以用 `\uhex` 指定 Unicode 码值,其中 `hex` 为码值的十六进制表示。所以,在 Unicode 编码环境下可以用字符组 `[\u4e00-\u9fff]` 匹配所有的中文字符。

```
/[\u4e00-\u9fa5]/.test("罚");    // => True
/[\u4e00-\u9fa5]/.test("a");    // => False
```

要注意的是,如果用点号 `.` 匹配多字节字符,在不同浏览器中的表现是不一样的,有时候可能匹配单个字节,有时候可能匹配整个多字节字符。在我的机器上,同一段 JavaScript 代码,Firefox 和 Chrome 中运行得到的结果就不一样。

#### Firefox

```
/^.$/.test("罚");    // => True
```

#### Chrome

```
/^.$/.test("罚");    // => False
```

## 12.2.3 字符组简记法

在 JavaScript 的正则表达式中, `\d`、`\D`、`\w`、`\W`、`\s`、`\S` 都采用 ASCII 匹配规则,即便在 Unicode 编码环境下仍是如此。也就是说, `\d` 等价于 `[0-9]`, `\w` 等价于 `[0-9a-zA-Z_]`, `\s` 无法匹配其他空白字符(比如中文的全角空格)。

```
//半角字符
/^\d$/ .test("1");    // => True
/^\w$/ .test("a");    // => True
/^\s$/ .test(" ");    // => True
//全角字符及中文字符
/^\d$/ .test(" 1");   // => False
/^\w$/ .test("我");   // => False
/^\s$/ .test(" ");    // => False
```

## 12.2.4 单词边界

JavaScript 的正则表达式中单词边界的规定与对 `\w` 的规定是一样的,也就是说, `\b` 匹配成功只有一种情况:一边必须保证 `\w` (等价于 `[0-9a-zA-Z_]`) 匹配成功,另一边必须保证 `\w` 匹配不成功。

```

//空字符串-英文字符
/a\b/.test("a");           // => True
//英文字符-半角标点
/a\b/.test("a,");          // => True
//英文字符-半角空格
/a\b/.test("a ");          // => True
//英文字符-中文字符
/a\b/.test("a 我");         // => True
//英文字符-全角空格
/a\b/.test("a ");          // => True
//中文字符-半角标点
/我\b/.test("我,");         // => False
//中文字符-半角空格
/我\b/.test("我 ");        // => False
//中文字符-中文字符
/我\b/.test("我和");        // => False
//中文字符-全角空格
/我\b/.test("我 ");        // => False

```

### 12.2.5 行起始/结束位置

在 JavaScript 的正则表达式中，能够用来定位字符串或行起始位置的只有 `^`（而没有 `\A`）。在默认情况下，`^` 只能匹配“整个字符串的开始位置”；如果指定使用多行模式（Multiline Mode），`^` 可以匹配字符串内部的文本行的开始位置。

```

/^1/.test("1\n2"); // => True
/^2/.test("1\n2"); // => False
/^2/m.test("1\n2");// => True

```

同样，能够定位字符串或行结束位置的只有 `$`（而没有 `\z` 和 `\Z`）。在默认情况下，`$` 只能匹配整个字符串的结束位置。如果指定使用多行模式，`$` 可以匹配字符串内部的文本行的结束位置。

```

/2$/.test("1\n2"); // => True
/1$/.test("1\n2"); // => False
/1$/m.test("1\n2");// => True

```

即便字符串末尾有换行符，`$` 也只匹配字符串的结束位置，而不会匹配末尾换行符之前的位置。所以如果进行正则表达式验证，可以放心地在表达式中使用 `^` 和 `$`。

```

/^\d{6}$/.test("123456\n"); // => False
/^\d{6}$/.test("123456");    // => True

```

### 12.2.6 环视

JavaScript 1.5 以后的版本（也就是说，除特别古老的 JavaScript 之外）都支持肯定顺序环视

`(?=...)`和否定顺序环视`(?!...)`，环视中的表达式并没有任何限制。

```
/c(?:=ab+)/.test("cab");           // => True
/c(?:=(ab+|cd))/test("cab");        // => True
/c(?:=(ab+|cd))/test("cabbb");      // => True
```

不过，JavaScript 的正则表达式中不支持肯定逆序环视`(?<=...)`和否定逆序环视`(?<!...)`，有时候，这确实非常不方便。

## 12.2.7 匹配模式

JavaScript 对匹配模式的支持非常有限，常见的匹配模式中只支持不区分大小写模式（`i`）和多行模式（`m`），且不支持在正则表达式中用`(?modifier)`的记法设定模式，只能在末尾的反斜线之后写上模式对应的字母，所以模式必须对整个表达式生效。

```
/a/.test("A");           // => False
/a/i.test("A");          // => True
/1$/m.test("1\n2");      // => True
```

也可以混用多个模式，可以把模式对应的字母连续列出来，指定不区分大小写模式和多行模式。

```
/^LINE2/im.test("line1\nline2"); // => True
```

如果使用字符串形式生成 `RegExp`，则必须提供第二个参数来设定模式。

```
(new RegExp("A", "i")).test("A");           // => True
(new RegExp("line1$", "m")).test("line1\nline2"); // => True
(new RegExp("^LINE2", "im")).test("line1\nline2"); // => True
```

因为JavaScript的正则表达式不支持单行模式，所以点号`.`在任何情况下都不能匹配换行符——在IE浏览器中，点号`.`不能匹配换行符`\n`，但可以匹配回车符`\r`；而在Firefox、Chrome、Safari等浏览器中，点号`.`既不能匹配`\n`，也不能匹配`\r`。

**IE**

```
/line1./.test("line1\rline2"); // => True
```

**Firefox Chrome Safari**

```
/line1./.test("line1\rline2"); // => False
```

如果你真的需要能够匹配任意字符，不妨使用`\s\S`或者`[\d\D]`，它们确实可以匹配包括换行符在内的任意字符。

```
//任何浏览器都成立
/line1[\s\S]/.test("line1\rline2"); // => True
//对换行符也成立
/line1[\s\S]/.test("line1\nline2"); // => True
```

JavaScript 中的另一种匹配模式是 **g**，它表示 **Global**（全局模式，对整个字符串执行），也就是说，正则表达式的操作不会进行停止于第一次匹配。在提取文本和替换时，这个模式很有用，下文会详细介绍。

## 12.2.8 捕获分组的引用

在 JavaScript 中，如果要在正则表达式内部引用捕获分组，则应当使用 `\num` 记法，其中 *num* 为对应捕获分组的编号。

```
/^([a-z])\1$/.test("ab");    // => False
/^([a-z])\1$/.test("aa");    // => True
```

如果要在替换的 *replacement* 字符串中引用捕获分组，则应当使用 `$num` 记法。

```
"2010-12-20".replace(/(\d{4})-(\d{2})-(\d{2})/, "$2/$3/$1");
12/20/2010
```

如果要在 *replacement* 字符串中表示 `$` 字符，则必须转义为 `$$`。

```
"the price is 12.99".replace(/\d+\.\d{0,2}/, "$$$0");
the price is $12.99
```

## 12.3 正则 API 简介

### 12.3.1 RegExp

JavaScript 中的正则表达式对象是 **RegExp**，生成它的方法主要有两种：使用 `/regex/` 分隔符表示法指定，或显式调用 `RegExp()` 构造函数。

```
//简便形式，使用正则文字，不需要考虑字符串转义
var pattern = /\d+/;
//指定了多行模式
var pattern = /\d+/m;
```

显式使用 `RegExp()` 构造函数有两种形式：直接以分隔符和正则文字表示，类似 `/regex/`；或者以字符串作为参数，这时候就必须考虑字符串转义。如果使用前一种形式，必须在结束分隔符之后指定匹配模式；如果使用后一种形式，只能通过第二个参数指定匹配模式。

```
//常见形式，使用正则文字
var pattern = new RegExp(/\d+/);
//以字符串形式给出，需要考虑字符串中的转义处理，\d 必须写成\\d
var pattern = new RegExp("\\d+");
//指定多行模式的两种方式
var pattern = new RegExp(/\d+/m);
var pattern = new RegExp("\\d+", "m");
```



第 6 章介绍过，一般情况下，编程语言中正则表达式对象和匹配对象是不同的，表达式对象不会包含与匹配有关的信息，但 JavaScript 中的 `RegExp` 是例外，它包含 5 个状态变量，如表 12-2 所示，其中的 `lastIndex` 尤其值得注意。

表 12-2 `RegExp` 中的变量

变量	说明
<code>source</code>	只读变量，保存文本形式的表达式
<code>global</code>	只读变量，标识正则表达式是否指定了全局模式（g）
<code>ignoreCase</code>	只读变量，标识正则表达式是否指定了不区分大小写模式（i）
<code>multiline</code>	只读变量，标识正则表达式是否指定了多行模式（m）
<code>lastIndex</code>	可写变量，如果使用了全局模式，这个变量保存的是在字符串中尝试下次匹配的偏移值，在 <code>RegExp.test()</code> 和 <code>RegExp.exec()</code> 中都会用到这个变量

如果使用 `RegExp` 的 `exec()` 和 `test()` 函数，并且设定了全局模式，正则表达式的匹配就会从 `lastIndex` 的位置开始，并且在每次匹配成功之后重新设定 `lastIndex`。这样，就可以在字符串中重复迭代，依次寻找各个匹配结果。但是，如果需要对不同的字符串调用同一个 `RegExp` 的 `exec()` 或 `test()` 方法，这个变量也有可能带来意料之外的匹配结果，所以应该记得更换字符串时，显式将 `RegExp` 的 `lastIndex` 设定为 0。

12.3.1.1 `RegExp.exec(string)`

这个方法用来在字符串中寻找匹配，如果成功，则返回表示匹配信息的数组，否则返回 `null`。在默认情况下，它返回第一次匹配的结果。

```
var pattern = /\d{4}-\d{2}-\d{2}/;
var str = "2010-12-20 2011-02-14";
if ((matchArray = pattern.exec(str))!= null) {
    document.write (matchArray[0] + " starts at " + matchArray.index + "<br />");
}
2010-12-20 starts at 0
```

上一节说过，`RegExp` 对象中存在 `lastIndex` 变量，它指定下次开始尝试匹配的位置。但是如果指定全局模式，同时把判断语句从 `if` 改为 `while`，那么每次调用 `RegExp.exec()` 时都会从字符串的起始位置开始尝试匹配，结果会造成无穷循环。

```
var pattern = /\d{4}-\d{2}-\d{2}/;
var str = "2010-12-20 2011-02-14";
while ((matchArray = pattern.exec(str))!= null) {
    document.write (matchArray[0] + " starts at " + matchArray.index + "<br />");
}
2010-12-20 starts at 0
2010-12-20 starts at 0
```

```
2010-12-20 starts at 0
.....
```

指定全局模式之后，RegExp 对象会每次匹配成功，都会把匹配的结束位置更到 `lastIndex`，下次调用时从 `lastIndex` 开始尝试匹配。

```
var pattern = /\d{4}-\d{2}-\d{2}/g;
var str = "2010-12-20 2011-02-14";
while ((matchArray = pattern.exec(str))!= null) {
    //注意匹配结束位置是使用 pattern.lastIndex 获得的
    document.write (matchArray[0] + " starts at " + matchArray.index + " and ends at "
    + pattern.lastIndex + "<br />");
}
2010-12-20 starts at 0 and ends at 10

2011-02-14 starts at 11 and ends at 21
```

你或许注意到了，上面的程序先把 RegExp 对象保存在变量 *pattern* 中，反复调用，而不是直接在 while 语句中写 `/\d{4}-\d{2}-\d{2}/g.exec(str)`。这是因为按照 ECMAScript 5 规范，除非复用 RegExp 对象，否则每次遇到 */regex/* 或 `new RegExp(regex)` 时，都会重新生成 RegExp 对象，这样就会造成死循环。

```
var str = "2010-12-20 2011-02-14";
while ((matchArray = /\d{4}-\d{2}-\d{2}/g.exec(str))!= null) {
    document.write (matchArray[0] + " starts at " + matchArray.index + "<br />");
}
2010-12-20 starts at 0
2010-12-20 starts at 0
2010-12-20 starts at 0
.....
```

在 `RegExp.exec(string)` 的匹配结果 `matchArray` 中，*index* 表示匹配文本在原文本中的开始位置，而 `matchArray` 数组中的各个元素，依次对应正则表达式中各个捕获分组匹配的文本。按照惯例，编号为 0 的分组对应整个正则表达式匹配的文本。在 JavaScript 中，只能获取各捕获分组匹配的文本，而不能获得各分组匹配文本的开始和结束位置。

```
var pattern = /(\d{4})-(\d{2})-(\d{2})/g;
var str = "2010-12-20 2011-02-14";
while ((matchArray = pattern.exec(str))!= null) {
    document.write (matchArray[0] + " starts at " + matchArray.index + " and ends at "
    + pattern.lastIndex + "<br />");
    document.write (matchArray[1] + "<br />");
    document.write (matchArray[2] + "<br />");
    document.write (matchArray[3] + "<br />");
}
```

```

}
2010-12-20 starts at 0 and ends at 10

2010

12

20

2011-02-14 starts at 11 and ends at 21

2011

02

14

```

### 12.3.1.2 RegExp.test(string)

此方法在本章开头已经介绍过，它用来测试正则表达式能否在字符串中找到匹配文本，返回布尔值。

```
/\d{4}-\d{2}-\d{2}/.test("2010-12-20")); // => True
```

在前端开发中，经常需要验证整个字符串是否能由正则表达式匹配。要完成这种功能，可以在表达式的两端分别加上匹配字符串起始和结束的`^`和`$`（注意：JavaScript 并不支持`\A`和`\Z`）。因为即便字符串末尾是换行符，`$`也会匹配真正的结束位置，所以可以放心使用。

```

/^\d{4}-\d{2}-\d{2}$/.test("2010-12-20")); // => True
/^\d{4}-\d{2}-\d{2}$/.test("a2010-12-20")); // => False
/^\d{4}-\d{2}-\d{2}$/.test("2010-12-20\n")); // => False

```

每个 `RegExp` 对象都包含状态变量 `lastIndex`。如果指定了全局模式，每次执行 `RegExp.test()` 时，都会从字符串中的 `lastIndex` 偏移值开始尝试匹配，所以用同一个 `RegExp` 多次验证字符串，必须记得每次调用之后，将 `lastIndex` 设定为 0，否则就可能出错了。当然，更好的办法是调用 `string.search(RegExp)`，判断返回值是否为 0。

```

//保证使用同一个 RegExp 对象，且指定 g 模式
var pattern = /\d{4}-\d{2}-\d{2}$/g;
pattern.test("2010-12-20")); // => True
pattern.test("2010-12-20")); // => False

```

`RegExp.test(string)` 还有一点奇妙之处：匹配成功之后，各捕获分组匹配的文本已经保存下来，用 `RegExp.$1`、`RegExp.$2`… 就可以获得，不过，“理应”保存整个表达式所匹配文本的 `RegExp.$0` 并不存在。

```

if (/^\d{4}-\d{2}-\d{2}$/.test("2010-12-20")) {
    document.write(RegExp.$2 + "/" + RegExp.$3 + "/" + RegExp.$1);
}
12/20/2010

```

应当注意，如果当前匹配不成功，`RegExp.$1` 等变量仍然可能有值，只不过是最近一次成功调用 `RegExp.test(string)` 时的值。

```

/^(\\d{4})-(\\d{2})-(\\d{2})$/.test("2010-12-20"); //匹配成功
/^(\\d{4})-(\\d{2})-(\\d{2})$/.test("abcd-ef-gh"); //匹配失败
document.write(RegExp.$2 + "/" + RegExp.$3 + "/" + RegExp.$1);
12/20/2010

```

## 12.3.2 String

JavaScript 中的一些正则操作（比如查找）既可以通过 `RegExp` 的方法实现，也可以通过 `String` 类的方法实现；另一些正则操作（比如切分和替换）是调用 `String` 类的方法实现的。

### 12.3.2.1 string.match(RegExp)

这个方法类似 `RegExp.exec(string)`，只是调换了 `RegExp` 和 `string` 的位置。两者的另一个区别是，无论是否指定全局模式，`Regexp.exec()` 总是返回单次的匹配结果，而 `string.match()` 在指定了全局模式的情况下会返回一个字符串数组，其中包含各次成功匹配的文本，但不包含任何其他信息。

```

var pattern = /\\d{4}-\\d{2}-\\d{2}/g;
var str = "2010-12-20 2011-02-14";
var matchArray = str.match(pattern);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}
2010-12-20

2011-02-14

```

需要注意的是，即便正则表达式中包含捕获分组，在 `string.match(RegExp)` 返回的数组中，也不会包含分组捕获的信息。

```

var pattern = /(\\d{4})-(\\d{2})-(\\d{2})/g;
var str = "2010-12-20 2011-02-14";
var matchArray = str.match(pattern);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}
2010-12-20

2011-02-14

```

### 12.3.2.2 string.search(RegExp)

这个方法用来寻找某个正则表达式在字符串中第一次匹配成功的位置，如果不成功，则返回-1。

```
"a1".search(/\d); // =>1
"ab".search(/\d/); // =>-1
```

这个方法只能找到“第一次”匹配的位置，即便设定了全局模式，结果也不会有任何变化。所以，如果需要用表达式多次验证字符串，调用这个方法判断结果是否等于0，是更好的办法。

```
"1234".search(/^\d+$/g) == 0; // => True
"5678".search(/^\d+$/g) == 0; // => True
```

相比之下，同样的 `RegExp` 对象，如果设定全局模式再调用 `RegExp.test(string)`，可能会因为 `RegExp` 中保存的 `lastIndex` 变量，得到错误的结果。

```
var pattern = /^\d+$/g;
pattern.test("1234"); // => True
//这时候 lastIndex 的值是 4
pattern.test("5678"); // => False
```

### 12.3.2.3 string.replace(RegExp, replacement)

这个方法用来进行正则表达式替换，将 *RegExp* 能匹配的文本替换为 *replacement*。在默认情况下，它只进行一次替换，如果设定了全局模式，则所有能匹配的文本都会替换。

```
"2010-12-20 2011-02-14".replace(/\d{4}-\d{2}-\d{2}/, "Date");
Date 2011-02-14

"2010-12-20 2011-02-14".replace(/\d{4}-\d{2}-\d{2}/g, "Date");
Date Date
```

如果要在 *replacement* 字符串中引用分组，则可以使用 *\$num*。

```
"2010-12-20".replace(/(\d{4})-(\d{2})-(\d{2})/, "$2/$3/$1");
12/20/2010
```

如果要在 *replacement* 字符串中表示 *\$* 字符，则必须使用 *\$\$* 转义。

```
"the price is 12.99".replace(/\d+\.\d{0,2}/, "$$$0");
the price is $12.99
```

在 *Replacement* 字符串中，还有一些特殊变量可以使用，表 12-3 总结了常用的特殊变量。

表 12-3 JavaScript 的正则表达式中常用的特殊变量

特殊变量	说明
\$num	表示对应捕获分组匹配的文本
\$\$	\$ 字符
\$`	正则表达式所匹配文本之前（左侧）的文本
\$'	正则表达式所匹配文本之后（右侧）的文本

12.3.2.4 string.replace(RegExp, function)

常见的文档中并没有说明这个方法，它其实可以算是上面方法的“变种”，只是把第二个参数从字符串变成了函数，这个函数的接收参数是一个字符串，返回参数也是一个字符串。这样一来，文本的处理就更加灵活了。下面的例子展示了如何将 t 开头的单词转换为大写。

```
document.write("one two three".replace(/\bt[a-zA-Z]+\b/g,
function(m){return m.toUpperCase()}));
one TWO THREE
```

12.3.2.5 string.split(RegExp)

这个方法使用一个正则表达式切分字符串，正则表达式是否使用了全局模式对结果没有影响。

```
var matchArray = "one two three".split(/\s+/);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}
one
two
three
```

也可以设置第二个参数 *limit*，手工指定返回数组的数目。需要注意的是，在其他语言中，如果 *limit* 小于实际能切分的次数，返回数组的最后那个元素一般是“正则表达式第 *limit*-1 次匹配右侧的所有文本”，但 JavaScript 中却不是这样，它仍然会进行尽可能多的切分，返回切分所得数组的第 *limit* 个元素。

```
//注意最后一个结果不是 two three
var matchArray = "one two three".split(/\s+/, 2);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}
one
two
```

表 12-4 总结了 *limit* 在各个取值区间对结果的影响（假定未指定 *limit* 时，最多可以切分 *n*-1 次，返回包含 *n* 个元素的数组）。

表 12-4 limit 各种取值的结果

取值	结果
limit < 0	等于未设定 limit 时，返回包含 <i>n</i> 个元素的数组
0 <= limit <= n	返回一个包含 <i>n</i> 个元素的数组，其中最后一个元素是切分 <i>n</i> 次时的第 limit 个元素（如果 limit 设定为 0，则返回空数组）
limit > n	等于未指定 limit 时

## 12.4 常用操作示例

### 12.4.1 验证

```
//必须在正则表达式首尾加上^和$
var regex = /^d{4}-d{2}-d{2}$/;
"2010-12-20".search(regex) == 0;      // => True

//如果使用 RegExp.test(), 一定不能设定全局模式
var regex = /^d{4}-d{2}-d{2}$/;
regex.test("2010-12-20");              // => True
regex.test("2011-02-14");              // => True
```

### 12.4.2 提取

#### 简单提取

```
var pattern = /d{4}-d{2}-d{2}/g;
var str = "2010-12-20 2011-02-14";
var matchArray = str.match(pattern);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}

2010-12-20

2011-02-14
```

#### 逐步提取

```
//如果要获取各捕获分组匹配的文本，只能使用这个方法，同时必须指定 g 模式
var pattern = /(\d{4})-(\d{2})-(\d{2})/g;
var str = "2010-12-20 2011-02-14";
while ((matchArray = pattern.exec(str))!= null) {
```

```

        document.write("date: " + matchArray[0]);
        document.write(", year: " + matchArray[1]);
        document.write(", month: " + matchArray[2]);
        document.write(", day: " + matchArray[3]);
        document.write("<br />");
    }
    date: 2010-12-20, year: 2010, month: 12, day: 20

    date: 2011-02-14, year: 2011, month: 02, day: 14

```

### 12.4.3 替换

#### 简单替换

```

var regex = /(\d{4})-(\d{2})-(\d{2})/;
var replacement = "$2/$3/$1";
document.write("2010-12-20".replace(regex, replacement));
12/20/2010

```

#### 在 replacement 中使用 \$

```

var regex = /\d+\.\d{0,2}/;
var replacement = "$$$0";
document.write("the price is 12.99".replace(regex, replacement));
the price is $12.99

```

### 12.4.4 切分

```

var regex = /\s+/;
var matchArray = "one\ntwo\t three".split(regex);
for(var i = 0; i < matchArray.length; i++) {
    document.write(matchArray[i] + "<br />");
}
one two three

```

## 12.5 关于 ActionScript

ActionScript 是用于 Adobe Flash 和 Adobe AIR 运行时环境的编程语言，它同样遵循 ECMA Script 标准，其中的正则表达式操作与 JavaScript 非常相似，所以在这里简要介绍 ActionScript 中的正则表达式。

### 12.5.1 RegExp

与 JavaScript 一样，ActionScript 中的正则对象同样是 **RegExp**，只是生成方式有细微区别，在 ActionScript 中，必须显式指定所生成对象的类型为 **RegExp**。



```
//通过正则文字生成
var pattern1:RegExp = /bob/;
//通过正则文字生成, 同时指定不区分大小写模式
var pattern1:RegExp = /bob/i;
//通过字符串生成
var pattern2:RegExp = new RegExp("bob");
//通过字符串生成, 同时指定不区分大小写模式
var pattern2:RegExp = new RegExp("bob", "i");
```

12.5.2 匹配规则

ActionScript 只支持 ASCII 匹配规则, 所以 `\w`、`\s`、`\d` 都无法匹配 ASCII 编码之外的字符。

```
//半角字符
/^\d$/.test("1"); // => True
/^\w$/.test("a"); // => True
/^\s$/.test(" "); // => True
//全角字符及中文字符
/^\d$/.test(" 1 "); // => False
/^\w$/.test("我"); // => False
/^\s$/.test(" "); // => False
```

12.5.3 匹配模式

相比 JavaScript, ActionScript 对匹配模式的支持要好很多, 它一共支持 5 种匹配模式, 如表 12-5 所示。

表 12-5 ActionScript 中的匹配模式

模式修饰符	说明
g	全局模式, 意义与 JavaScript 的相同, 用于提取多个结果
i	不区分大小写模式
s	单行模式
m	多行模式
x	注释模式

同样, ActionScript 中的 RegExp 对象也会保留 `lastIndex` 信息。所以, 如果你希望对某个字符串重复应用同一个正则表达式, 依次查找, 请记得指定全局模式。

12.5.4 正则 API

ActionScript 中的正则 API 与 JavaScript 中的几乎完全一致, 为避免重复, 这里不做重复说明。

## 第 13 章 PHP

PHP 内建了三套各自独立的正则引擎，分别是 `preg`、`ereg`、`mb_ereg`。不过在日常使用最多的是 `preg`（`preg` 表示 PCRE Regex，PCRE 表示 Perl Compatible Regular Expression，常见的正则表达式写法都是兼容 Perl 的），本章的讲解也以 `preg` 为准<sup>1</sup>。

PHP 现在有两大分支，分别是 PHP 4.x 和 PHP 5.x，它们各成一系，没有绝对的高低之分。不过关于正则表达式的部分，PHP 4.x 和 PHP 5.x 基本是没有差别的，遇到有差异的部分会详细说明。

### 13.1 预备知识

PHP 的正则表达式是以字符串文字指定的，不过又不同于普通的字符串文字，如果用一段字符串文字表示正则表达式，文字的首位必须加上分隔符（*delimiter*），通常使用的分隔符是斜线/，在 4.0.4 以上的版本中，也可以采用 `()`、`{}`、`<>`、`[]` 这些成对的符号。不过 `()`、`{}`、`[]` 等符号都是正则表达式中常用的元字符，如果用作分隔符，用作元字符时就也必须转义，非常烦琐，而在正则表达式中没有特殊含义，所以一般推荐用/做分隔符，本章中也是如此。

```
$regex = '/regex/'; //这是正确的写法
$regex = '{regex}'; //这也是正确的写法
```

---

<sup>1</sup> 如果不确定机器上是否有 PCRE，可以按照下面的方法测试，这里附上我的测试结果。

```
yusheng@yurii_ubuntu:~$ pcretest -C
PCRE version 7.8 2008-09-05
Compiled with
    UTF-8 support
    Unicode properties support
    Newline sequence is LF
    \R matches all Unicode newlines
    Internal link size = 2
    POSIX malloc threshold = 10
    Default match limit = 10000000
    Default recursion depth limit = 10000000
    Match recursion uses stack
```

```
$regex = 'regex'; //这是常犯的错误，正则表达式没错，但忘了加分隔符
```

你可能注意到了，上面的字符串文字是用单引号标注的。在 PHP 中，字符串既可以用双引号标注，也可以用单引号标注。两者的主要区别在于，双引号字符串可以插值，而单引号字符串不能；另外，双引号字符串会处理字符串转义，而单引号字符串不会。

```
$regex = "[a-z]";
echo "double-quoted\t$regex";
echo "<br />";
echo 'single-quoted\t$regex';

double-quoted [a-z]
single-quoted\t$regex
```

如果不需要插值，使用双引号字符串的处理效率就不如单引号字符串，再加上双引号字符串用作正则表达式时还需要处理字符串转义，所以本章表示正则表达式时，都使用单引号字符串，直接写出正则文字，而字符串文字则使用双引号字符串，方便识别 `\r`、`\n` 之类常见的转义序列。

下面介绍正则表达式匹配的函数 `preg_match($pattern, $subject)`。它接收两个参数，其中 *pattern* 是正则表达式，*subject* 是需要匹配的文本；返回的是匹配发生的次数——如果能找到匹配则立刻停止，返回 1，否则返回 0；如果匹配过程中发生错误，则返回 `false`。需要注意的是，它并不检查整个 *subject* 能否由 *pattern* 匹配。

```
preg_match('/\d/', "1"); // => True
preg_match('/\d/', "a1b"); // => False
```

注：因为 PHP 可以把 `int` 类型作为布尔判断的条件，0 等价于 `false`，其他值等价于 `true`，下面的代码里，为了形象起见，如果只关心 `preg_match()` 的返回值是否为 0，使用 `true` 和 `false` 来标识。

在正则表达式两端添加 `\A` 和 `\Z`，就可以判断整个 *subject* 能否由 *pattern* 匹配。

```
preg_match('/\A\d\Z/', "1"); // => True
preg_match('/\A\d\Z/', "a1b"); // => False
```

在有些 PHP 代码中，使用双引号字符串来表示正则表达式，某些时候，情况看起来有点奇怪。

```
preg_match("/\d/", "1"); // => True
preg_match("/\d/", "1"); // => True
```

尽管 PHP 的正则表达式用字符串文字给出，但它与常见的字符串不完全一样，主要是转义规则有所区别——如果某个转义序列可以由字符串识别，则对其进行转义处理；否则，将整个转义序列“原封不动”地保存下来。在处理正则表达式时，“原封不动”保存下来的文字，会作为正则文字，只进行正则转义处理（详见第 93 页）。

这样现象可能很难理解：`/\d/`和 `/\d/`看起来不同，却表示同样的正则表达式。所以，我推荐的做法是（理由见第 6 章）：正则表达式尽量使用单引号字符串，比如 `'/\d/'`或`'/\A/'`；如果

是在要使用双引号字符串，尽量写为"/\\d/"或"/\\A/"，因为大多数语言中的都可以这么写。

关于 mb\_string

用过 PHP 的人大都知道 PHP 中的 mb\_string，它用来处理多字节字符构成的字符串 (multibyte-string)，不会将多字节字符拆散为多个单字节字符。比如用 mb\_strpos()来计算偏移值，得到偏移值的单位就是字符，而不是字节。

```
strpos("位置", "置", 0);           //3
mb_strpos("位置", "置", 0, "utf-8"); //1
```

相应的，mb\_string 中也提供了正则表达式函数，比如 mb\_ereg\_match()、mb\_ereg\_replace()、mb\_regex\_encoding()等，其功能类似对应的 preg 函数，具体情况请参考 <http://www.php.net/manual/en/ref.mbstring.php>。

按照文档说明，在设定 mb\_regex\_encoding()之后，可以正确处理非 Unicode 编码的正则表达式，但实际使用中往往容易出现各种奇怪的问题（例如上面提到的[收发]的错误匹配，即便设定了 GBK 编码，也会出现），所以一般不推荐这样使用正则表达式。

13.2 正则功能详解

13.2.1 列表

表 13-1 中列出了 PHP 中的正则功能。

表 13-1 PHP 中的正则功能

功能	记法	说明
字符组 ①1	[...] [^...]	支持 Unicode
POSIX 字符组 ①15	[[:digit:]]	只匹配 ASCII 字符
Unicode 属性 ①121	\p{...}	支持
字符组简记法 ①115	\d \D \w \W \s \S	采用 ASCII 匹配规则
行起始位置 ①61	^ \A	支持
行结束位置 ①61	\$ \z \Z	支持
单词边界 ①117	\b \B	采用 ASCII 匹配规则
顺序环视 ①68	(?=...) (?!...)	支持
逆序环视 ①68	(?<=...) (?<!...)	不支持无法确定长度范围的表达式
匹配模式 ①81	i m s x	支持

( 续表 )

功能	记法	说明
模式作用范围 89	<code>(?modifier) (?-modifier)</code>	支持
纯文本模式 99	<code>\Q...\E</code>	支持
捕获分组及引用 44	<code>(...) \num \$num</code>	支持
命名分组 53	<code>(?&lt;name&gt;...) \k&lt;name&gt;</code>	支持
非捕获分组 54	<code>(?:...)</code>	支持
多选结构 39	<code>(... ...)</code>	支持
匹配优先量词 18	<code>? * +</code>	支持
忽略优先量词 26	<code>?? ?* ?+ {n,m}?</code>	支持

13.2.2 字符组

PHP 中的正则引擎可以正确识别 ASCII 编码，但无法正确识别 GBK 编码，因此无法避免 GBK 编码环境中的“错误匹配”问题（8111）。

```
preg_match('/[收发]/', "罚"); // => True
preg_match('/[收发]/', "发"); // => True
```

解决的办法是采用 Unicode 编码环境（以 Unicode 编码保存源文件），并在正则表达式末尾的分隔符之后加上模式修饰符 `u`，它指定 Unicode 模式，也就是按 UTF-8 编码解释正则表达式。

```
preg_match('/[收发]/u', "罚"); // => False
```

只有在指定 Unicode 模式之后，点号才能准确匹配多字节字符。

```
# GBK 编码环境
preg_match('/\A.\z/u', "罚"); // => False
# Unicode 编码环境
preg_match('/\A.\z/u', "罚"); // => True
```

PHP 的字符组支持使用 `\x{hex}` 通过码值表示字符，其中的 `hex` 是 4 位十六进制数，所以可以用 `[\x{4e00}-\x{9fff}]` 来匹配汉字字符，但是此时一定不能忘记 `u` 修饰符。

```
preg_match('/[\x{4e00}-\x{9fff}]/u', "我"); // => True
preg_match('/[\x{4e00}-\x{9fff}]/u', "a"); // => False
preg_match('/[\x{4e00}-\x{9fff}]/', "我"); //错误
```

在 PHP 中也可以直接使用 POSIX 字符组，所有的 POSIX 字符组都只能匹配 ASCII 字符，即便指定了 Unicode 模式也是这样。

```
preg_match('/[[:digit:]]/', "0"); // => True
preg_match('/[[:xdigit:]]/', "f"); // => True
```

### 13.2.3 Unicode 属性

从 PHP 5.1.0 开始, PHP 支持使用 Unicode 属性, 方法是 `\p{N}` 和 `\P{N}`, 其中 `N` 为 Unicode Property (☞123), 而且此时一定不能忘记 `u` 修饰符<sup>1</sup>。

`\p{N}` 可以匹配“任何一个具备指定 Unicode Property 的字符”, 而 `\P{N}` 匹配“任何一个不具备指定 Unicode Property 的字符”, 以表示大写字母的 `Lu` 为例。

```
preg_match('/\p{Lu}/u', "A"); // => True
preg_match('/\p{Lu}/u', "a"); // => False
```

```
preg_match('/\P{Lu}/u', "A"); // => False
preg_match('/\P{Lu}/u', "a"); // => True
```

`\P{N}` 也可以等价表示为 `\p{^N}`:

```
preg_match('/\p{^Lu}/u', "A"); // => False
preg_match('/\p{^Lu}/u', "a"); // => True
```

因为 PHP 中可以使用 Unicode Script (☞123), 所以可以用 `\p{Han}` 来匹配中文字符。

```
preg_match('/\p{Han}/u', "我"); // => True
```

### 13.2.4 字符组简记法

在 PHP 中, `\d`、`\D`、`\w`、`\W`、`\s`、`\S` 都采用字符组简记法的 ASCII 匹配规则, 即便指定了 Unicode 模式, 也是如此。

```
//半角字符
preg_match('/\d/', "0"); // => True
preg_match('/\w/', "a"); // => True
preg_match('/\s/', " "); // => True
//全角字符
preg_match('/\d/u', "0 "); // => False
preg_match('/\w/u', " a "); // => False
preg_match('/\s/u', " "); // => False
```

### 13.2.5 单词边界

在 PHP 中, `\w` 采用 ASCII 匹配规则, 所以在 `\b` 能匹配位置的某一侧必须出现 `[0-9a-zA-Z_]`, 另一侧不能出现 `[0-9a-zA-Z_]`。

```
//空字符串-英文字符
preg_match('/a\b/', "a"); // => True
```

---

<sup>1</sup> 请参考 <http://www.php.net/manual/en/regexp.reference.unicode.php>。

```
//英文字符-半角标点
preg_match('/a\b/', "a,"); // => True
//英文字符-半角空格
preg_match('/a\b/', "a "); // => True
//英文字符-中文字符
preg_match('/a\b/', "a 我"); // => True
//英文字符-全角空格
preg_match('/a\b/', "a "); // => True
//中文字符-半角标点
preg_match('/我\b/', "我,"); // => False
//中文字符-半角空格
preg_match('/我\b/', "我 "); // => False
//中文字符-中文字符
preg_match('/我\b/', "我和"); // => False
//中文字符-全角空格
preg_match('/我\b/', "我 "); // => False
```

前面讲过，PHP 处理双引号字符串时，会将字符串文字中无法识别的转义序列保留下来，所以“`\\d`”和“`\\d`”是没有差别的。值得注意的是“`\\b`”和“`\\b`”也是没有差别的，因为 PHP 的字符串转义并不能识别转义序列`\\b`。

### 13.2.6 行起始/结束位置

在 PHP 的正则表达式中，`^`匹配的是“行开始的位置”，在默认情况下，它只能匹配“整个字符串的开始位置”；如果指定使用多行模式（Multiline Mode），`^`可以匹配字符串内部的文本行的开始位置；而`^A`无论在什么情况下都匹配“整个字符串的开始位置”。

```
preg_match('/^1/', "1\n2\n"); // => True
preg_match('/^2/', "1\n2\n"); // => False
preg_match('/^2/m', "1\n2\n"); // => True

preg_match('/^A1/', "1\n2\n"); // => True
preg_match('/^A2/', "1\n2\n"); // => False
preg_match('/^A2/m', "1\n2\n"); // => False
```

在默认情况下，PHP 中的`$`、`\Z`、`\z`匹配的是整个字符串的结束位置（如果结束位置有换行符，则`$`和`\Z`匹配这个换行符之前的位置）；使用多行模式会影响到`$`的匹配，这时候它可以匹配文本内部行的结束位置；如果要进行严格准确的验证，则应当使用`\z`。

```
preg_match('/1$/', "1\n"); // => True
preg_match('/1$/', "1\n2"); // => False
preg_match('/1$/m', "1\n2"); // => True

preg_match('/1\Z/', "1\n"); // => True
preg_match('/1\Z/', "1\n2"); // => False
```

```
preg_match('/1\Z/m', "1\n2");    // => False

preg_match('/1\z/', "1");        // => True
preg_match('/1\z/', "1\n");      // => False
preg_match('/1\z/m', "1\n");     // => False
```

13.2.7 环视

PHP 中的肯定顺序环视( `?(=...)` )和否定顺序环视( `?(!...)` )内可以嵌套任意形式的子表达式。

```
preg_match('/c(?(=ab+)/', "cab");    // => True
preg_match('/c(?(=(ab+|cd))/', "ccd"); // => True
preg_match('/c(?(=(ab+|cd))/', "cabbbb"); // => True
```

PHP 的逆序环视中能使用的表达式，必须匹配固定长度的文本，否则会报错。如果表达式中出现了多选结构，而各个分支能匹配的文本长度不一致，那么这个多选结构必须位于表达式的最顶层，而且只能出现竖线|，类似 `bull|donkey`，而不能出现括号，比如 `(bull|donkey)`。

```
preg_match('/ab(?(<=ab)\/', "ab");    // => True
preg_match('/cd(?(<=ab|cd)\/', "cd"); // => True

//长度不固定，报错
preg_match('/cd(?(<=(dogs?|cats?))\/', "cd");
//多选结构各分支匹配文本长度不一致，且两端出现了括号，报错
preg_match('/cd(?(<=(bull|donkey))\/', "cd");
```

13.2.8 匹配模式

表 13-2 中列出了 PHP 中常用的匹配模式。

表 13-2 PHP 中常用的匹配模式

匹配模式	修饰符	解释
不区分大小写模式	i	不区分 ASCII 字符的大小写
注释模式	x	容许正则表达式中出现注释，表达式中的空白字符，以及# 开始到行末的文本，都视为注释
多行模式	m	^和\$不仅仅匹配字符串的起始/结束位置，还可以匹配字符串内部文本行的起始/结束位置
单行模式	s	容许点号.匹配任何字符，包括换行符
Unicode 模式	u	UTF-8 模式，不会把正则表达式中多字节字符拆散为字节，而是将其视为单个字符

PHP 中指定匹配模式的方式有两种：在结尾分隔符之后加上模式对应的字母；或者在正则表达式的开头用 `(?modifier)` 的方式表示。注意：在 PHP 中没有通过预定义常量指定模式的方法。



```
preg_match('/a/', "A");    // => False
preg_match('/a/i', "A");    // => True
preg_match('/(?i)a/', "A"); // => True
```

其中 Unicode 模式有些特殊,它只能在末尾分隔符之后使用,不能使用 `(?modifier)` 的记法。

```
preg_match('/\A.\z/u', "字");    // => True
preg_match('/(?u)\A.\z/', "字");  //编译错误
```

PHP 也支持用 `(?-modifier)` 结束某个模式的作用范围。

```
preg_match('/a(?i)b(?-i)b/', "abb");    // => True
preg_match('/a(?i)b(?-i)b/', "aBb");    // => True
preg_match('/a(?i)b(?-i)b/', "aBB");    // => False
```

### 13.2.9 纯文本模式

在 `\Q` 和 `\E` 之内,所有的元字符和特殊结构都失去特殊意义,只能匹配它们对应的字符本身。

```
preg_match('/\Q.\E/', "A");    // => False
preg_match('/\Q.\E/', ".");    // => True
```

### 13.2.10 捕获分组的引用

在 PHP 中,如果要在正则表达式内部引用捕获分组,则应当使用 `\num` 记法,其中 `num` 为对应捕获分组的编号。

```
preg_match('/([a-z])\1/', "ab");    // => False
preg_match('/([a-z])\1/', "aa");    // => True
```

如果要在替换时引用捕获分组,则应当使用 `$num` 记法,不过更好的办法是使用 `${num}` 记法,这样可以避免二义性(☞50)。

```
preg_replace('/(\\d{4})-(\\d{2})-(\\d{2})/', "$2/$3/$1", "2010-12-20");
12/20/2010

preg_replace('/(\\d{4})-(\\d{2})-(\\d{2})/', "${2}/${3}/${1}", "2010-12-20");
12/20/2010
```

如果希望在 *replacement* 字符串中表示 `$` 字符,则必须使用 `\\$` 转义。

```
preg_replace('/\d+\.\d{0,2}/', '\\${0}', "the price is 12.99");
the price is $12.99
```

PHP 支持命名分组,分组的记法是 `(?P<name>...)`,在表达式中引用的记法为 `(?P=name)` (在 5.2.2 版本以后可以使用 `\k<name>` 或者 `\k' name'`,在 5.2.4 版本之后可以使用 `\k{ name}` 和 `\g{ name}`),在替换中,却只能通过数字编号来引用。

```
preg_match('/(?P<char>[a-z])(?P=char)/', "ab"); // => False
preg_match('/(?P<char>[a-z])(?P=char)/', "aa"); // => True
```

需要指出的是，如果进行正则表达式替换，*replacement* 中不能使用命名分组，只能使用数字编号分组。

## 13.3 正则 API 简介

### 13.3.1 PREG 常量说明

preg 系列函数有可能用到预定义的常量，所以先介绍这些常量。

#### 13.3.1.1 PREG\_PATTERN\_ORDER

仅用于 `preg_match_all()`，在匹配结果的 `matches` 数组中，`matches[0]` 保存每次匹配成功时捕获的全部文本（也就是编号为 0 的捕获分组匹配的文本）；`matches[1]` 保存每次匹配成功时编号为 1 的捕获分组匹配的文本；依此类推。

#### 13.3.1.2 PREG\_SET\_ORDER

仅用于 `preg_match_all()`，在匹配结果的 `matches` 数组中，`matches[0]` 保存每一次匹配成功时，整个表达式以及各捕获分组匹配的文本；`matches[1]` 保存第 2 次匹配成功时，整个表达式以及各捕获分组匹配的文本；依此类推。

#### 13.3.1.3 PREG\_SPLIT\_OFFSET\_CAPTURE

在 PHP 4.3.0 以上版本中可用，仅用于 `preg_split()` 中。如果设定了这个常量，每次匹配之后，会返回后面的字符串的偏移值。请注意，返回值是一个数组，其中每个元素都是一个数组，包含偏移值为 0 的匹配的字符串，以及它在文本中的偏移值，从 1 开始。

#### 13.3.1.4 PREG\_OFFSET\_CAPTURE

在 PHP 4.3.0 以上版本中可用，其功能类似 `PREG_SPLIT_OFFSET_CAPTURE`，但它用于 `preg_match()` 和 `preg_matchall()`，如果此时指定 `PREG_SPLIT_OFFSET_CAPTURE`，则会报错。

#### 13.3.1.5 PREG\_SPLIT\_NO\_EMPTY

在 `preg_split()` 中使用它，返回结果中不会包含空字符串。

#### 13.3.1.6 PREG\_SPLIT\_DELIMI\_CAPTURE

在 PHP 4.0.5 以上版本中可用，在 `preg_split()` 中使用它，则会捕获正则表达式中的捕获括号所匹配的文本。

下面几个常量都是 `preg_last_error()` 返回的。

### 13.3.1.7 PREG\_NO\_ERROR

在 PHP 5.2.0 以上版本中可用，如果没有错误，则返回它。

### 13.3.1.8 PREG\_INTERNAL\_ERROR

在 PHP 5.2.0 以上版本中可用，如果 PCRE 发生了内部错误，则返回它。

### 13.3.1.9 PREG\_BACKTRACK\_LIMIT\_ERROR

在 PHP 5.2.1 以上版本中可用，如果尝试匹配时超过了回溯限制，则返回它。

### 13.3.1.10 PREG\_RECURSION\_LIMIT\_ERROR

在 PHP 5.2.0 以上版本中可用，如果匹配时超过了递归限制，则返回它。

### 13.3.1.11 PREG\_BAD\_UTF8\_ERROR

在 PHP 5.2.0 以上版本中可用，如果在 UTF-8 模式下正则表达式中出现了不合法的 UTF-8 数据，则返回它。

### 13.3.1.12 PREG\_BAD\_UTF8\_OFFSET\_ERROR

在 PHP 5.3.0 以上版本中可用，如果在 UTF-8 模式下起始偏移值对应的不是有效的 UTF-8 码值，则返回它。

### 13.3.1.13 PCRE\_VERSION

在 PHP 5.2.4 以上版本中可用，返回 PCRE 版本号和发布日期。

## 13.3.2 preg\_quote

```
string preg_quote ( string $str [, string $delimiter = NULL ] )
```

这个函数用来“包装”字符串 *str*。如果需要把这个字符串用作正则表达式，又要保证其中元字符都没有特殊含义（比如处理用户输入的内容），则可以使用它。

```
echo preg_quote('^.$');  
\\^\\.\\$
```

虽然在 PHP 中使用正则表达式时必须在首尾加上分隔符（*delimiter*），但分隔符可以使用多种字符。所以，`preg_quote()` 中的 *delimiter* 是可选参数；尽管一般用斜线/做分隔符，但在默认情况下，*str* 中的斜线并不会转义，除非显式指定使用/作为分隔符。

```
echo preg_quote('/');  
/  
  
echo preg_quote('/', '/');  
\/
```

如果要由某个字符串生成正则表达式，同时保证其中不包含任何元字符，可以用下面的办法（别忘了加上分隔符）。

```
$delimiter = '/';
$regex = $delimiter.preg_quote('.*/', $delimiter).$delimiter;
preg_match($regex, '.*/*'); // => True
```

也可以直接在字符串两端加上 `\Q` 和 `\E`。

```
$delimiter = '/';
$regex = $delimiter.'\Q'.origin.'\E'.$delimiter;
preg_match($regex, '.*/*'); // => True
```

### 13.3.3 preg\_grep

```
array preg_grep ( string $pattern, array $input [, int $flags = 0])
```

这个函数用来“筛选”数组，*input* 数组中所有能由正则表达式 *pattern* 匹配的字符串会被筛选出来，重新组成一个数组，作为结果返回<sup>1</sup>。

```
preg_grep('/^\d$/', array("1", "a", "2", "b"));
Array{ [0]=>"1", [2]=>"2" }
```

如果可选参数 *flags* 设定为常量 `PREG_GREP_INVERT`，则进行“反向筛选”，也就是说，*input* 中所有 *pattern* 不能匹配的字符串都被筛选出来，重新组成一个数组，作为结果返回。

```
preg_grep('/^\d$/', array("1", "a", "2", "b"), PREG_GREP_INVERT);
Array{ [0]=>"a", [3]=>"b" }
```

*flags* 参数在 PHP 4.2.0 以后才加入。

### 13.3.4 preg\_match

```
int preg_match ( string $pattern, string $subject [, array &$matches
[, int $flags = 0 [, int $offset = 0 ]]] )
```

这是很常用的函数，它在 *subject* 中查找 *pattern* 的匹配，返回匹配的次数。但是，它会在第一次匹配成功之后停下来，直接返回 1；如果始终不能匹配成功，则返回 0；如果发生错误，则返回 false。

```
preg_match('/\d/', "0123"); // => True
preg_match('/\d/', "abcd"); // => False
```

---

<sup>1</sup> 在本章中，为方便展示返回数组的内容，数组使用 `var_dump()` 函数展开的形式来表现。

根据 `preg_match()` 的返回值, 可以方便地进行数据验证, 但是注意在正则表达式的首尾加上 `\A` 和 `\Z`, 确保它验证的是整个字符串。下面以验证 6 到 8 位数字的字符串为例。

```
preg_match('/\A\d{6,8}\Z/', "012345");           // => True
preg_match('/\A\d{6,8}\Z/', "0123456789");       // => False
preg_match('/\A\d{6,8}\Z/', "a012345");          // => False
```

可选参数 *matches* 用来保存匹配结果, 如果匹配成功, 它包含一个元素, 即首次匹配的结果, 否则数组为空。

```
$matches = array();
preg_match('/\d/', "0123", $matches);
$matches: Array{ [0]=>"0" }

$matches = array();
preg_match('/\d/', "abcd", $matches);
$matches: Array{ }
```

如果 *pattern* 中出现了捕获分组, 则在 *matches* 数组中, 与捕获分组编号对应的元素, 保存该捕获分组匹配的文本。

```
$matches = array();
preg_match('/(\d)(\d)/', "0123", $matches);
$matches: Array{ [0]=>"01", [1]=>"0", [2]=>"1" }
```

可选参数 *flags* 控制匹配结果中是否保存匹配发生的位置, 如果设置为常量 `PREG_OFFSET_CAPTURE`, 则保存匹配发生的位置, 否则可以设置为 0。

```
$matches = array();
preg_match('/\d/', "0123", $matches, PREG_OFFSET_CAPTURE);
$matches: Array{ [0]=> Array{ [0]=>"0", [1]=>0 } }

$matches = array();
preg_match('/\d/', "0123", $matches, 0);
$matches: Array{ [0]=>"0" }
```

应该注意的是, 位置的单位是字节, 而不是字符。

```
$matches = array();
preg_match('/\d/', "字 0123", $matches, PREG_OFFSET_CAPTURE);
$matches: Array{ [0]=> Array{ [0]=>"0", [1]=>4 } }
```

可选参数 *offset* 控制匹配从 *subject* 的哪个位置开始。

```
$matches = array();
preg_match('/\d/', "0123", $matches, PREG_OFFSET_CAPTURE, 2);
$matches: Array{ [0]=> Array{ [0]=>"2", [1]=>2 } }
```

注意：这里位置的单位也是字节，而不是字符。

```
$matches = array();
preg_match('/\d/', "0 和 123", $matches, PREG_OFFSET_CAPTURE, 4);
$matches: Array{ [0]=> Array{[0]=>"1", [1]=>4} }
```

*flags* 参数和 `PREG_OFFSET_CAPTURE` 都是 PHP 4.3.0 之后加入的，*offset* 参数则是 PHP 4.3.3 之后加入的。

### 13.3.5 preg\_match\_all

```
int preg_match_all (string $pattern, string $subject, array
&$matches[, int $flags = PREG_PATTERN_ORDER[, int $offset = 0 ]])
```

这个函数与 `preg_match()` 相似，不同的是，它可以一次性找出正则表达式 *pattern* 在字符串 *subject* 中所有的匹配，返回匹配的次數；如果始终无法匹配，则返回 0；如果出错，则返回 false。在 `preg_match_all()` 中，同样也可以用 *matches* 数组来保存匹配结果。

```
$matches = array();
preg_match_all('/\d{4}-\d{2}-\d{2}/', "2010-12-20 2011-02-14", $matches);
$matches: Array{ [0]=>"2010-12-20", [1]=>"2011-02-14" }
```

如果 *pattern* 中出现了捕获型括号，则 *matches* 的每个元素都是数组，它保存对应编号的捕获分组在历次匹配中匹配到的文本。

```
$matches = array();
preg_match_all('/(\d{4})-(\d{2})-(\d{2})/', "2010-12-20 2011-02-14", $matches);
$matches: Array {
    [0]=> array(2)
        { [0]=>"2010-12-20", [1]=>"2011-02-14" },
    [1]=> array(2)
        { [0]=>"2010", [1]=>"2011" },
    [2]=> array(2)
        { [0]=>"12", [1]=>"02" },
    [3]=> array(2)
        { [0]=>"20", [1]=>"14" }
}
```

可以看到，*matches* 数组的第 0 个元素对应两次匹配时整个 *pattern* 匹配的文本（也就是第 0 个捕获分组匹配的文本），第 1 个元素对应两次匹配时 *pattern* 中第 1 个捕获分组匹配的文本，第 2 个元素对应两次匹配时 *pattern* 中第 2 个捕获分组匹配的文本。这是默认的方式，*flags* 显式设定为 `PREG_PATTERN_ORDER` 也是这样。

如果你在其他语言中使用过正则表达式，可能会觉得这样非常别扭，因为一般的逻辑中，匹配结果集中的每个元素对应到“每次匹配”的结果，而不是“每个分组匹配”的结果。如果希望

采用更符合惯例的逻辑，需要将 *flags* 设定为 `PREG_SET_ORDER`。

```
$matches = array();
preg_match_all('/(\d{4})-(\d{2})-(\d{2})/', "2010-12-20 2011-02-14", $matches,
PREG_SET_ORDER);
$matches: Array {
    [0]=> array(4)
        { [0]=>"2010-12-20", [1]=>"2010", [2]=>"12", [3]=>"20" },
    [1]=> array(4)
        { [0]=>"2011-02-14", [1]=>"2011", [2]=>"02", [3]=>"14" }
}

$matches = array();
preg_match_all('/(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})/', "2010-12-20 2011-02-14",
$matches, PREG_SET_ORDER);
```

```
$matches: Array {
    [0]=> Array {
        [0]=>"2010-12-20", ["year"]=>"2010", [1]=>"2010",
        ["month"]=> "12", [2]=>"12", ["day"]=>"20", [3]=>"20" },
    [1]=> Array {
        [0]=>"2011-02-14", ["year"]=>"2011", [1]=>"2011",
        ["month"]=>"02", [2]=>"02", ["day"]=>"14", [3]=>"14" }
}
```

*flags* 参数也可以设定为 `PREG_OFFSET_CAPTURE`，这样 *matches* 的结果相比 `PREG_PATTERN_ORDER`，就多出了匹配发生的位置信息（同样，单位是字节）。

```
$matches = array();
preg_match_all('/(\d{4})-(\d{2})-(\d{2})/', "2010-12-20 2011-02-14", $matches,
PREG_OFFSET_CAPTURE);
$matches: Array {
    [0]=> array(2)
        { [0]=>"2010-12-20", [1]=>"2011-02-14" },
    [1]=> array(2)
        { [0]=>"2010", [1]=>"2011" },
    [2]=> array(2)
        { [0]=>"12", [1]=>"02" },
    [3]=> array(2)
        { [0]=>"20", [1]=>"14" }
}
```

`PREG_OFFSET_CAPTURE` 都是 PHP 4.3.0 之后加入的，而 *offset* 参数则是 PHP 4.3.3 之后加入的，如果你使用的 PHP 版本不够高，可能无法使用这两项功能。

### 13.3.6 preg\_last\_error

这个函数返回最近一次调用 `preg` 系列方法时出现的错误，如果没有出错，则返回 `PREG_NO_ERROR`，其他错误可以参考“`PREG` 常量说明”。

### 13.3.7 preg\_replace

`mixed preg_replace ( mixed $pattern , mixed $replacement , mixed $subject [, int $limit = -1 [, int &$count ] ] )`

这个函数用来进行正则表达式替换，将 *subject* 中的 *pattern* 替换为 *replacement*。*pattern* 用于匹配需替换文本的正则表达式，可以是单个字符串，也可以是字符串数组；*replacement* 表示希望替换成的文本，可以在其中引用匹配的相关信息。如果要引用分组，推荐使用 `${num}` 的记法，这种记法没有二义性，不会引起混淆。

```
preg_replace('/(\d{4})-(\d{2})-(\d{2})/', '${2}/${3}/${1}', "2010-12-20 2011-02-14");
12/20/2010 02/14/2011
```

如果要在 *replacement* 中使用 `$` 符号，则必须使用 `\$` 转义。

```
preg_replace('/\d+\.\d{0,2}/', '\${$0}', "the price is 12.99");
the price is $12.99
```

*subject* 表示要进行替换操作的对象，它可以是单个字符串，也可以是字符串数组。如果是字符串，则返回的也是字符串；如果是数组，则返回的是数组，同时数组中的每个元素都会进行替换操作。

另外还有两个可选参数，其中 *limit* 在 PHP 4.0.2 以上版本中可用，表示 *pattern* 在 *subject* 中最多能匹配的次数，默认值为 -1，表示不限次数；*count* 在 PHP 5.1.0 以上版本中可用，如果传入了这个参数，匹配完成后会将它赋值为匹配实际发生的次数。

### 13.3.8 preg\_replace\_callback

`mixed preg_replace_callback ( mixed $pattern , callback $callback , mixed $subject [, int $limit = -1 [, int &$count ] ] )`

这个函数的功能类似 `preg_replace()`，不同的是，它不会简单把 *pattern* 能匹配的文本替换成其他文本，而是通过 *callback* 函数来处理，这样大大增加了灵活性，*callback* 接收的参数是一个 *match* 数组（而不是字符串），返回字符串。下面的代码借助这个函数，把所有单词统一为首



字母大写格式。

```
function capitalize($matches) {
    //正则表达式已经用两个捕获分组分开第一个字母和之后的字母
    return strtoupper($matches[1]).strtolower($matches[2]);
}
preg_replace_callback( '/\b([a-z])([a-z+)]\b/i', "capitalize", "one TWO three");
One TWO THREE
```

如果觉得这样很麻烦，也可以用 `create_function()` 定义匿名函数<sup>1</sup>。

```
preg_replace_callback( '/\b([a-z])([a-z+)]\b/i', create_function('$matches',
'return strtoupper($matches[1]).strtolower($matches[2]); '), "one TWO three");
```

在 PHP 5.3.0 以上版本中还可以直接定义匿名函数<sup>2</sup>。

```
preg_replace_callback( '/\b([a-z])([a-z+)]\b/i', function($matches) {return
strtoupper($matches[1]).strtolower($matches[2]);} , "one TWO three");
```

### 13.3.9 preg\_filter

`mixed preg_filter (mixed $pattern, mixed $replacement, mixed $subject [, int $limit = -1 [, int &$amp;count ]])`

这个函数类似 `preg_replace()`，唯一的区别在于，它返回的是确实发生过查找-替换的字符串：如果 *subject* 是字符串，且发生了查找-替换，则返回替换之后的字符串，否则返回 `NULL`；如果 *subject* 是数组，且其中有元素发生了查找-替换，则返回这些发生了查找-替换的元素构成的数组，否则返回空数组。

```
$textArray = array("2010-12-20", "other text", "2011-02-14");
preg_filter('/(\d{4})-(\d{2})-(\d{2})/', "$2/$3/$1", $textArray);

Array( [0]=>"12/20/2010", [1]=>"02/14/2011" )
```

PHP 5.2.0 之后提供了另一个函数 `filter_var()`，它通过预定义的过滤器来判断字符串的值，进行常用的过滤，而不需要显式用到正则表达式，如果判断成功则返回字符串，否则返回 `false`。常见的过滤器包括：

**FILTER\_VALIDATE\_EMAIL**：验证是否为 E-mail

**FILTER\_VALIDATE\_FLOAT**：验证是否为浮点型

---

<sup>1</sup> 可以参考 <http://www.php.net/manual/en/function.create-function.php>。

<sup>2</sup> 可以参考 <http://www.php.net/manual/en/functions.anonymous.php>。

**FILTER\_VALIDATE\_INT**: 验证是否为整型

**FILTER\_VALIDATE\_IP**: 验证是否是合法 IP

**FILTER\_VALIDATE\_URL**: 验证是否是合法 URL

以下是示例代码:

```
var_dump(filter_var("one@example....com", FILTER_VALIDATE_EMAIL));
false

var_dump(filter_var("one@example.com", FILTER_VALIDATE_EMAIL));
one@example.com
```

### 13.3.10 preg\_split

`array preg_split(string $pattern, string $subject [, int $limit = -1 [, int $flags = 0 ]])`

这个函数用来切分字符串, 返回切分所得的字符串数组。*pattern* 匹配的文本作为切分的分隔, 而 *subject* 则是切分的字符串。

```
preg_split('/\s+/', "one\ttwo\n three");
Array{ [0]=>"one", [1]=>"two", [2]=>"three" }
```

*limit* 是可选参数, 它设定 *pattern* 匹配的次数。如果设定了 *limit*, 那么返回数组包含 *limit* 个元素, 也就是说, 最多进行 *limit*-1 次切分。如果 *limit* 设为 -1、0、null, 都表示“不设定上限”, 也就是“切分尽可能多”。

```
preg_split('/\s+/', "one\ttwo\n three", -1);
Array{ [0]=>"one", [1]=>"two", [2]=>"three" }
```

```
preg_split('/\s+/', "one\ttwo\n three", 2);
Array{ [0]=>"one", [1]=>"two three" }
```

*flags* 是可选参数, 用于控制返回的数据, 还可以选择常量 `PREG_SPLIT_NO_EMPTY`, 它要求返回数组忽略空字符串。

```
preg_split('/', '/', "one,two,,three", -1);
Array{ [0]=>"one", [1]=>"two", [2]=>"", [3]=>"three" }
```

```
preg_split('/', '/', "one,two,,three", -1, PREG_SPLIT_NO_EMPTY);
Array{ [0]=>"one", [1]=>"two", [2]=>"three" }
```

有可能引起混淆的是 *limit* 和 `PREG_SPLIT_NO_EMPTY` 同时作用的情况, *limit* 设定的并不是“返回数组的长度”, 而是“*pattern* 匹配的次数”, 所以并不会“先全部切分, 再过滤数组”。

```
preg_split('/./', "one,two,,three", 2, PREG_SPLIT_NO_EMPTY);
Array{ [0]=>"one", [1]=>"two", [2]=>","three" }
```

## 13.4 常见的正则操作举例

### 13.4.1 验证

```
//记得要在正则表达式首尾添加\A 和\z
$regex = '/\A\d{4}-\d{2}-\d{2}\z/';
preg_match($regex, "2010-12-20"); // => True
```

### 13.4.2 提取

逐个提取

```
$matches = array();
$offset = 0;
$regex = '/\d{4}-\d{2}-\d{2}/';
$string = "2010-12-20 2011-02-14";
while (preg_match($regex, $string, $matches, PREG_OFFSET_CAPTURE, $offset)!=0 ) {
    echo($matches[0][0]);
    echo "<br />";
    $offset = $matches[0][1] + strlen($matches[0][0]);
}
2010-12-20
2011-02-14

$matches = array();
$offset = 0;
$regex = '/(\d{4})-(\d{2})-(\d{2})/';
$string = "2010-12-20 2011-02-14";
while (preg_match($regex, $string, $matches, PREG_OFFSET_CAPTURE, $offset)!=0 ) {
    echo "date: ".$matches[0][0].", year: ".$matches[1][0].", month: ".$matches[2][0].", day:
".$matches[3][0];
    echo "<br />";
    $offset = $matches[0][1] + strlen($matches[0][0]);
}
date: 2010-12-20, year: 2010, month: 12, day: 20
date: 2011-02-14, year: 2011, month: 02, day: 14
```

### 一次性提取

```
//为符合下关，均设定 PREG_SET_ORDER 以每次匹配为单位
$matches = array();
$regex = '/\d{4}-\d{2}-\d{2}/';
$string = "2010-12-20 2011-02-14";
preg_match_all($regex, $string, $matches, PREG_SET_ORDER) ;
foreach ($matches as $match) {
    echo($match[0]);
}

2011-03-17

2010-12-24

$matches = array();
$offset = 0;
$regex = '/(\d{4})-(\d{2})-(\d{2})/';
$string = "2010-12-20 2011-02-14";
preg_match_all($regex, $string, $matches, PREG_SET_ORDER) ;
foreach ($matches as $match) {
    echo("date: ".$match[0].", year: ".$match[1].", month: ".$match[2].", day:
".$match[3]);
    echo("<br />");
}

date: 2010-12-20, year: 2010, month: 12, day: 20

date: 2011-02-14, year: 2011, month: 02, day: 14
```

### 13.4.3 替换

#### 简单替换

```
$regex = '/(\d{4})-(\d{2})-(\d{2})/';
$string = "2010-12-20 2011-02-14";
$replacement = '${2}/${3}/${1}';
preg_replace($regex, $replacement, $string);
12/20/2010 02/14/2011
```

#### 在 replacement 中使用\$

```
$regex = '/\d+\.\d{0,2}/';
$replacement = '\${$0}';
$string = "the price is 12.99";
preg_replace($regex, $replacement, $string);
the price is $12.99
```

### 13.4.4 切分

```
$regex = '/-/';  
$string = "2010-12-20";  
preg_split($regex, $string);  
Array { [0]=>"2010", [1]=>"12", [2]=>"20" }
```

## 第 14 章 Python

Python 对正则表达式的支持较为完备，它的 API 很丰富，也支持 Unicode 字符串（Python 2 的字符串并没有默认使用 Unicode 编码，使用稍微麻烦），但不支持 Unicode 属性（☞121）。因为 Python 2 和 Python 3 在正则表达式的规定上有较大区别，在区别较大的环节会分版本讲解<sup>1</sup>。

### 14.1 预备知识

Python 内置了正则表达式的 package，在使用正则表达式之前，必须首先导入对应 package：

```
import re
```

为讲解方便，先介绍 Python 中的 `re.search(pattern, string)` 函数。如果 `pattern` 可以匹配 `string` 的某个子串，则返回一个 `MatchObject` 对象，否则返回 `None`。通过判断是否返回 `None`，就可以知道 `pattern` 能否匹配 `string` 中的子串。如果在正则表达式两端加上 `\A` 和 `\Z`（注意是 `\Z` 而不是 `\z`，Python 中没有 `\z`，`\Z` 等价于其他语言中的 `\Z`），则可以判断整个字符串是否能由正则表达式匹配，进行正则表达式验证。

```
re.search("\\d{6}", "123456") != None      # => True
re.search("\\d{6}", "12345678") != None    # => True
re.search("\\d{6}", "a123456") != None     # => True
re.search("\\A\\d{6}\\Z", "123456") != None # => True
re.search("\\A\\d{6}\\Z", "12345678") != None # => False
re.search("\\A\\d{6}\\Z", "a123456") != None # => False
```

Python 的字符串比较复杂，值得花点篇幅介绍一下：Python 中常用的字符串有单引号字符串和双引号字符串，两者并没有区别，`"\\d{6}"` 与 `'\\d{6}'` 是完全等价的。如果使用这两种字符串来表示正则表达式，都必须经过字符串转义与正则转义（详见第 6 章），也就是说，正则表达式 `\\d` 应当表现为字符串文字 `\\d`。

但是 Python 也提供了另一种“原生字符串（Raw String）”，其方式是在字符串之前添加字符 `r`，这样就会忽略字符串转义，只需要提供正则文字即可。这样，正则表达式 `\\d` 就可以写作 `r"\\d{6}"` 或 `r'\\d'`。这种形式更适合正则表达式，所以本章中的正则表达式都使用这种形式。

---

<sup>1</sup> Python 2 和 Python 3 的另一大差别在于 `print` 在 Python 3 中成为函数，所以 `print text` 之类的写法不再可行，而应当改为 `print(text)`。考虑到目前 Python 2 和 Python 3 的份额，本章仍然采用了 Python 2 的 `print` 写法。

```
re.search(r"\A\d{6}\Z", "123456") != None      # => True
re.search(r"\A\d{6}\Z", "12345678") != None    # => False
re.search(r"\A\d{6}\Z", "a123456") != None     # => False
```

你有可能在某些代码中看到正则表达式 `\d` 直接使用字符串文字 `\d`，而没有使用原生字符串，但并没有出错。这是因为 Python 在处理字符串转义时，如果字符串文字中的某个转义序列无法识别，会原封不动地“保留下来”。`\d` 并不是可以识别的字符串转义序列，所以会被保留下来，作为正则文字。因此 `\d` 也可以直接写成字符串 `\d`，同样 `\A` 和 `\Z` 也是如此（☞93）。

```
re.search("\A\d{6}\Z", "123456") != None      # => True
re.search("\A\d{6}\Z", "12345678") != None    # => False
re.search("\A\d{6}\Z", "a123456") != None     # => False
```

但我并不推荐使用这种做法，因为表达式 `\b` 不能写成 `"\b"`——因为在字符串文字中 `\b` 是可以识别的字符串转义序列，表示退格符（backspace），所以 `"\bcat\b"` 生成的正则表达式不能匹配“单词 cat”，而是匹配“两个退格符之间的 cat”。

```
re.search("\bcat\b", "cat") != None           # => False
re.search("\bcat\b", "cat") != None           # => True
re.search(r"\bcat\b", "cat") != None           # => True
```

另外，Python 2 的字符串并非默认采用 Unicode 编码，所以如果要用到 Unicode 功能（并非 Unicode 属性）<sup>1</sup>，则必须在字符串之前添加字符 `u`，比如 `u"[你我他]"`。所以，如果字符串中出现了多字节字符，最好使用 `u` 属性，写成 `u"你"`；如果正则表达式中出现了多字节字符，最好是同时使用 `u` 和 `r`。不过，注意：只能写作 `ur"[你我他]"`，不能写作 `ru"[你我他]"`。而在 Python 3 中，因为字符串默认采用 Unicode 编码，所以只需用 `r` 标记，比如 `r"[你我他]"`。

## 14.2 正则功能详解

### 14.2.1 列表

表 14-1 中列出了 Python 中的正则功能。

表 14-1 Python 中的正则功能

功能	记法	说明
字符组 ☞1	<code>[...]</code> <code>[^...]</code>	支持 Unicode
POSIX 字符组 ☞15	<code>[[:digit:]]</code>	只匹配 ASCII 字符

<sup>1</sup> Python 2.x 中并没有“Unicode 字符串”，而只有“Unicode 对象”，所谓“Unicode 字符串”也只是被当做 `str` 的“Unicode 对象”而已。而 `str` 对象只是一个字符数组，里面存储什么内容，使用什么编码，都需要用户自己负责。所以，如果在 Windows 的 Python IDLE（而不是命令行提示符）下使用 Python，即便指定了 `u`，也生成不了真正的 Unicode 字符串，测试本章中的例子可能会出错。

(续表)

功能	记法	说明
Unicode 属性 121	<code>\p{...}</code>	不支持
字符组简记法 115	<code>\d \D \w \W \s \S</code>	2.x 中默认采用 ASCII 匹配规则，但可设定 Unicode 模式变为 Unicode 匹配规则；3.x 中默认使用 Unicode 模式，但可设定 ASCII 模式变为 ASCII 匹配规则
行起始位置 61	<code>^ \A</code>	支持
行结束位置 61	<code>\$ \z \Z</code>	不支持 <code>\z</code>
单词边界 117	<code>\b \B</code>	2.x 中默认采用 ASCII 匹配规则，但可设定 Unicode 模式变为 Unicode 匹配规则；3.x 中默认使用 Unicode 模式，但可设定 ASCII 模式变为 ASCII 匹配规则
顺序环视 68	<code>(?=...)</code> <code>(?!...)</code>	支持
逆序环视 68	<code>(?&lt;=...)</code> <code>(?&lt;!...)</code>	不支持变长表达式
匹配模式 81	<code>i m s x a u</code>	支持
模式作用范围 89	<code>(?modifier)</code> <code>(?-modifier)</code>	只支持 <code>(?modifier)</code>
纯文本模式 99	<code>\Q...\E</code>	不支持
捕获分组及引用 44	<code>(...)</code> <code>\num \$num</code>	支持，同时支持命名分组
命名分组 53	<code>(?P&lt;name&gt;...)</code> <code>(?P=name)</code>	支持
非捕获分组 54	<code>(?:...)</code>	支持
多选结构 39	<code>(... ...)</code>	支持
匹配优先量词 18	<code>?</code> <code>*</code> <code>+</code>	支持
忽略优先量词 26	<code>??</code> <code>?*</code> <code>?+</code> <code>{n,m}?</code>	支持

14.2.2 字符组

14.2.2.1 Python 2

除非显式指定，否则 Python 2 的字符串默认都采用 ASCII 编码，而且也不支持 GBK 编码环境下使用正则表达式，就可能遇到错误匹配的问题（111）。

```
re.search(r"\A.\Z", "发") != None      # => False
re.search(r"[收发]", "罚") != None      # => True
```



解决的办法是使用 Unicode 编码环节。请注意，不管是用作正则表达式的字符串，还是被匹配的字符串，只要出现了中文字符，都必须指定 `u`。

```
re.search(r"\A.\Z", u"发") != None      # => True
re.search(ur"[收发]", u"罚") != None    # => False
```

### 14.2.2.2 Python 3

Python 3 中的字符串默认采用 Unicode 编码，所以不存在上面的问题。

### 14.2.2.3 关于编码

如果不清楚自己使用的 Python 是采用什么编码，可以运行下面的代码得到。

```
import sys
sys.getdefaultencoding()
```

如果显示 ASCII，则会遇到 GBK 编码的错误；如果显示为 UTF-8，则不会。

有时，Python 可能会提示程序文件中出现的 UTF-8 字符无法识别，这时应该在.py 文件的最开头加上这样一行：

```
# encoding=utf-8
```

在 Mac OS 和 Linux/Unix 下，Python 会自动从环境变量中获得编码（通过设定 `LANG=en_US.UTF-8` 或 `LANG=zh_CN.UTF-8`）。Windows 默认使用的并不是 Unicode 编码，所以，如果要在 Windows 下使用 Python，又要使用 UTF-8 编码，除了直接使用 Python 3，Python 2 必须在.py 文件的顶部加上这行。

## 14.2.3 Unicode 属性

Python 2 和 Python 3 都不支持 Unicode 属性，也就是说，所有规定的 Unicode 属性都无法在 Python 中使用。

## 14.2.4 字符组简记法

### 14.2.4.1 Python 2

在 Python 2 中，在默认情况下，`\d`、`\D`、`\w`、`\W`、`\s`、`\S` 都采用字符组简记法的 ASCII 匹配规则（☞115）。也就是说，`\d` 等价于 `[0-9]`，`\w` 等价于 `[0-9a-zA-Z_]`，`\s` 无法匹配其他空白字符（比如中文的全角空格）。

```
#半角字符
re.search(r"\d", "1") != None      # => True
re.search(r"\w", "a") != None      # => True
re.search(r"\s", " ") != None      # => True
```

```
#全角字符及中文字符
re.search(r"\d", u" 1 ") != None # => False
re.search(r"\w", u"我") != None # => False
re.search(r"\s", u" ") != None # => False
```

但是如果启用 Unicode 模式（最简单的做法是在表达式最前面加上(?u)），则会启用字符组简记法的 Unicode 匹配规则。

```
#全角字符及中文字符
re.search(r"(?u)\d", " 1 ") != None # => True
re.search(r"(?u)\w", "我") != None # => True
re.search(r"(?u)\s", " ") != None # => True
```

到现在为止，出现了两个概念：“Unicode 字符串”和“Unicode 模式”，它们比较容易混淆，表 14-2 有助于澄清这两个概念。

表 14-2 Python 中的 Unicode 字符串和 Unicode 模式

名称	Unicode 字符串	Unicode 模式
使用对象	字符串	正则表达式
作用	把汉字作为单个字符对待，而不是割裂为多个字节	扩展字符组简记法中“数字”、“空白”、“单词字符”的意义
形式	在字符串前面加上 u，比如 u"我"	在表达式开头加上(?u)，比如 u"(?u)\w"
应用场景	字符串中出现了汉字	需要 <code>\d</code> 、 <code>\D</code> 、 <code>\w</code> 、 <code>\W</code> 等简记法识别中文字符
可否连用	可以，但必须记为 ur，不能记为 ru	

14.2.4.2 Python 3

在 Python 3 中，各字符组简记法默认都采用字符组简记法的 Unicode 匹配规则（☞115）。也就是说，在默认情况下，`\d` 能匹配 Unicode 的数字字符，`\w` 能匹配 Unicode 的单词字符，`\s` 能匹配 Unicode 的空白字符。

```
#半角字符
re.search(r"\d", "1") != None # => True
re.search(r"\w", "a") != None # => True
re.search(r"\s", " ") != None # => True

#全角字符及中文字符
re.search(r"\d", " 1 ") != None # => True
re.search(r"\w", "我") != None # => True
re.search(r"\s", " ") != None # => True
```

如果要“还原”为 ASCII 匹配规则，则必须显式设定 ASCII 模式（最简单的做法是在表达式

最前面加上(?a))。

```
#全角字符及中文字符
re.search(r"(?a)\d", "1") != None    # => False
re.search(r"(?a)\w", "我") != None    # => False
re.search(r"(?a)\s", " ") != None     # => False
```

在 Python 2 和 Python 3 中, 默认的字符组简记法匹配规则完全相反, 使用时必须注意。比如在 Python 2 中完全可以用 `^\d+$` 验证 ASCII 字符数字的字符串, 但 Python 3 中则不能这样做。

## 14.2.5 单词边界

### 14.2.5.1 Python 2

Python 2 的正则表达式中单词边界的规定与对 `\w` 的规定是一样的。也就是说, `\b` 匹配成功只有一种情况: 一边必须保证 `\w` (等价于 `[0-9a-zA-Z_]`) 匹配成功, 另一边必须保证 `\w` 匹配不成功。

```
#空字符串-英文字符
re.search(r"a\b", "a") != None        # => True
#英文字符-半角标点
re.search(r"a\b", "a,") != None       # => True
#英文字符-半角空格
re.search(r"a\b", "a ") != None       # => True
#英文字符-中文字符
re.search(r"a\b", "a 我") != None     # => True
#英文字符-全角空格
re.search(r"a\b", "a ") != None       # => True
#中文字符-半角标点
re.search(r"我\b", "我,") != None     # => False
#中文字符-半角空格
re.search(r"我\b", "我 ") != None     # => False
#中文字符-中文字符
re.search(r"我\b", "我和") != None    # => False
#中文字符-全角空格
re.search(r"我\b", "我 ") != None     # => False
```

如果显式指定了 Unicode 匹配模式, 则 `\w` 能匹配的字符包括中文字符, 所以 `\b` 的匹配也相应变化。

```
#英文字符-中文字符
re.search(r"(?u)a\b", "a 我") != None # => False
#中文字符-半角标点
re.search(r"(?u)我\b", "我,") != None  # => True
#中文字符-半角空格
```

```
re.search(r"(\u)我\b", "我 ") != None    # => True
#中文字符-全角空格
re.search(r"(\u)我\b", "我 ") != None    # => True
```

### 14.2.5.2 Python 3

因为 Python 3 默认采用 Unicode 匹配规则，默认情况下，`\w` 认定的单词字符就包括其他语言中的字符(包括中文字符)，等价于 Python 2 显式指定 Unicode 模式的情况。如果希望使用 ASCII 匹配规则，则应当显式指定 ASCII 模式，它等价于 Python 2 中的默认情况。

## 14.2.6 行起始/结束位置

`^`匹配的是“行开始的位置”，在默认情况下，它只能匹配“整个字符串的开始位置”，如果指定使用多行模式 (Multiline Mode)，`^`可以匹配字符串内部的文本行的开始位置；而`\A`无论在什么情况下都匹配“整个字符串的开始位置”。

```
re.search(r"^\n", "1\n2\n") != None      # => True
re.search(r"^\2", "1\n2\n") != None      # => False
re.search(r"^(?m)^\2", "1\n2\n") != None # => True

re.search(r"^\A1", "1\n2\n") != None     # => True
re.search(r"^\A2", "1\n2\n") != None     # => False
re.search(r"^(?m)\A2", "1\n2\n") != None # => False
```

在默认情况下，Python 中的`$`和`\Z`匹配的是整个字符串的结束位置(如果结束位置有换行符，则`$`匹配这个换行符之前的位置)，使用多行模式会影响`$`的匹配，这时候它可以匹配文本内部行的结束位置。所以如果要进行严格准确的验证，则应当使用`\Z`，它等价于其他语言中的`\Z`。

```
re.search(r"1$/", "1\n") != None         # => True
re.search(r"1$/", "1\n") != None         # => False
re.search(r"^(?m)1$ ", "1\n2") != None   # => True

re.search(r"1\Z", "1") != None           # => True
re.search(r"1\Z", "1\n") != None         # => False
re.search(r"^(?m)1\Z", "1\n2") != None   # => False
```

## 14.2.7 环视

Python 中的肯定顺序环视`(?=...)`和否定顺序环视`(?!...)`内可以使用任意形式的子表达式。

```
re.search(r"c(?:=ab+)", "cab")           # => True
re.search(r"c(?:=(ab+|cd))", "ccd")      # => True
re.search(r"c(?:=(ab+|cd))", "cabbbb")   # => True
```

Python 中逆序环视则多了很多限制：逆序环视`(?<=...)`和`(?!<...)`中，表达式能匹配的文本长度必须是固定的，否则就会报错。也就是说，逆序环视中的表达式不能使用`+`、`*`、`?`之类的量

词；如果出现多选结构，各个多选分支能匹配的文本长度也必须是相同的。

```
re.search(r"ab(?:<=ab)", "ab")           # => True
re.search(r"cd(?:<=ab?)", "cd");           #报错，ab?能匹配的文本长度不确定
re.search(r"cd+(?:<=(abcd|cd))", "cd");    #报错，abcd 和 cd 匹配的文本长度不同
```

14.2.8 匹配模式

Python 的正则表达式支持所有常见的模式，另外还提供了 Unicode 模式和 ASCII 模式，详见表 14-3。

表 14-3 Python 中的匹配模式

常量	修饰符	解释
re.I re.IGNORECASE	i	不区分大小写模式
re.X re.VERBOSE	x	注释模式
re.M re.MULTILINE	m	多行模式
re.S re.DOTALL	s	单行模式
re.U re.UNICODE	u	Unicode 模式
re.A re.ASCII	a	ASCII 模式 <sup>1</sup>

注：re.A 和 re.ASCII 模式只存在于 Python 3 中。

匹配模式可以在表达式中以 `(?...)` 指定，也可以将常量作为参数指定。

```
re.search("a", "A") != None           # => False
re.search("(?i)a", "A") != None       # => True
re.search("a", "A", re.IGNORECASE) != None # =>True
```

需要注意的是，Python 中的模式永远是对整个正则表达式生效的，无论 `(?...)` 出现在哪里。

```
re.search("a(?i)b", "AB") != None # => True
```

Python 不支持用 `(?-...)` 停用模式。

14.2.9 捕获分组的引用

在 Python 中，如果要在正则表达式内部引用捕获分组，则应当使用 `\num` 记法，其中 *num*

<sup>1</sup> re.A 和 re.ASCII 模式只有在 Python 3 中才有。

为对应捕获分组的编号。

```
re.search(r"([a-z])\1", "ab") != None    # => False
re.search(r"([a-z])\1", "aa") != None    # => True
```

如果要在替换时引用捕获分组，则仍然使用`\num` 记法。

```
re.sub(r"(\d{4})-(\d{2})-(\d{2})", r"\2/\3/\1", "2010-12-20")
12/20/2010
```

更好的办法是使用`\g<num>`，这样更加清晰，避免了二义性。

```
re.sub(r"(\d{4})-(\d{2})-(\d{2})", r"\g<2>/\g<3>/\g<1>", "2010-12-20")
12/20/2010
```

如果使用了命名分组，在表达式中应当使用 `(?P=name)` 来引用，在替换时应当使用 `\g<name>` 来引用。

```
re.search(r"(?P<char>[a-z])(?P=char)", "ab") != None    # => False
re.search(r"(?P<char>[a-z])(?P=char)", "aa") != None    # => True

regex = r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})";
replacement = r"\g<month>/\g<day>/\g<year>";
print re.sub(regex, replacement, "2010-12-20")

12/20/2010
```

## 14.3 正则 API 简介

### 14.3.1 RegexpObject

Python 中的正则表达式对象是 `RegexpObject`，它可以由正则表达式编译而来，常用的方法是 `re.compile(pattern)`。

```
regexObject = re.compile("ab+")
```

如果要指定匹配模式，则可以在表达式中使用模式修饰符，也可以使用预定义常量。

```
regexObject = re.compile("(?i)ab+")
regexObject = re.compile("ab+", re.I)
```

如果要同时指定多个匹配模式，则可以在表达式中连写模式修饰符，也可以用 `|` 运算符连接预定义常量。

```
regexObject = re.compile("(?is)"ab+")
regexObject = re.compile("ab+", re.I|re.S)
```

在 Python 的文档中，大多数函数是类似 `re.search(pattern, string)` 的形式，其中的 *pattern* 通常是正则表达式的字符串，这可以视为“静态方法”。其实也可以将 `RegexObject` 对象作为 *pattern* 参数传入，或者直接对 `RegexObject` 对象调用相同名字的成员方法，结果相同，但是这两种做法能显式缓存 `RegexObject` 对象。下面的代码分别使用了三种办法，结果相同。

```
re.search(r"\d{6}", "123456") != None      # => True

sixDigitRegex = re.compile(r"\d{6}")
sixDigitRegex.search("123456") != None     # => True

sixDigitRegex = re.compile(r"\d{6}")
re.search(sixDigitRegex, "123456") != None  # => True
```

`re.compile()` 还有个不为许多人所知的功能，就是用来观察某个正则表达式的详细信息，方法是指定第二个参数为 `re.DEBUG`。如果你遇到复杂的表达式，或者不确定某个表达式的意义，可以通过它来观察。下面给出两个简单的例子，其中的缩进表示了表达式各结构的层级关系，而字符本身则显示为其码值的十进制表示（比如 `a` 的码值是十进制的 97）。

```
re.compile("ab|cd", re.DEBUG)
branch
    literal 97
    literal 98
or
    literal 99
    literal 100

re.compile("(ab|[cde])+", re.DEBUG)
max_repeat 1 65535
subpattern 1
    branch
        literal 97
        literal 98
    or
        in
            literal 99
            literal 100
            literal 101
```

### 14.3.2 re.compile(regex[, flags])

这个方法用于显式“编译”正则表达式对象。如果需要多次使用同一个正则表达式，那么，每次使用编译好的正则表达式对象，比每次临时编译表达式要快得多。`re.compile()`的实际使用在上一节已经看到过，这里不再重复。

需要指出的是，Python 会自动缓存最近一次在 `re.match()`、`re.search()` 和 `re.compile()` 中指定的正则表达式。如果用到的正则表达式并不多，可以使用字符串形式的正则表达式，而不必显式调用 `re.compile()` 获得 `RegexObject` 对象。也正因为如此，本章的大部分例子都没有显式调用 `re.compile()`，而是直接以字符串方式给出正则表达式。在实际编程中，可以根据具体情况决定是否显式调用 `re.compile()`。

### 14.3.3 re.search(pattern, string[, flags])

这个方法用来测试正则表达式 *pattern* 能否在 *string* 中找到匹配，*flags* 是可选参数，表示匹配模式（下面各个方法中的 *flags* 参数都具有同样的功能，不再赘述）。如果能找到匹配，则返回 `MatchObject`（详见下文），否则返回 `None`。

```
re.search(r"\d", "ab") #匹配不成功
re.search(r"\d", "ab1") #得到一个 MatchObject
<_sre.SRE_Match object at 0x0000000001D55E00>
```

这个函数可以用来进行数据验证，在正则表达式两端加上 `\A` 和 `\Z`（Python 不支持 `\z`，`\Z` 等价于其他语言中的 `\z`），并判断返回值是否为 `None`。

```
re.search(r"\A\d{6,10}\Z", "1234567ab") != None # => False
re.search(r"\A\d{6,10}\Z", "1234567") != None # => True
```

### 14.3.4 MatchObject

调用 `re.search()` 或者 `re.match()`，如果匹配成功，就返回一个 `MatchObject` 对象。实际上，与正则表达式匹配有关的许多方法都会返回 `MatchObject` 对象。

`MatchObject` 提供了若干方法和属性，通过它们可以获得匹配的信息，表 14-4 列出了常用的方法和属性。

表 14-4 MatchObject 常用的方法和属性

方法	描述	注释
<code>start(<i>n</i>)</code>	返回编号为 <i>n</i> 的捕获分组匹配的起始位置，如果对应分组不存在或未匹配，则返回 -1	如果没有设定 <i>n</i> ，则返回整个表达式匹配的起始位置，此时 <i>n</i> 等价于 0
<code>pos</code>	返回整个表达式匹配的起始位置	



( 续表 )

方法	描述	注释
end( <i>n</i> )	返回编号为 <i>n</i> 的捕获分组匹配的结束位置，如果对应分组不存在或未匹配，则返回-1	如果没有设定 <i>n</i> ，则返回整个表达式匹配的结束位置，此时 <i>n</i> 等价于 0
endpos	返回整个表达式匹配的结束位置	等价于 end(0)
group( <i>n</i> )	返回编号为 <i>n</i> 的捕获分组匹配的文本	如果没有设定 <i>n</i> ，则返回整个表达式匹配的文本，此时 <i>n</i> 等价于 0
groups()	返回各捕获分组匹配的文本构成的 tuple	如果表达式中不存在捕获分组，则返回空 tuple
span( <i>n</i> )	返回编号为 <i>n</i> 的捕获分组匹配的开始-结束位置	如果没有设定 <i>n</i> ，则返回整个表达式匹配的开始-结束位置，此时 <i>n</i> 等价于 0
re	返回匹配时所使用的正则表达式对象	
lastindex	返回匹配成功的编号最大的捕获分组的编号	如果没有，则返回 None
string	返回用来尝试用正则表达式来匹配的字符串	
expand( <i>str</i> )	返回一个字符串，可以在其中引用捕获分组匹配的文本	

下面的代码集中示范了其中主要的方法。

```
obj = re.search(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-20")
print "%s starts at %d and ends at %d" % (obj.group(), obj.pos, obj.endpos)
for i in range(1, obj.lastindex+1) :
    print "%s starts at %d and ends at %d" % (obj.group(i), obj.start(i), obj.end(i))
print obj.expand(r"year:\1 month:\2 day:\3")
2010-12-20 starts at 0 and ends at 10
2010 starts at 0 and ends at 4
12 starts at 5 and ends at 7
20 starts at 8 and ends at 10
year:2010 month:12 day:20
```

14.3.5 re.match(pattern, string[, flags])

re.match()和 re.search()非常相似，参数相同，返回值也相同，非常容易混淆，唯一的区别在于：如果 re.match()匹配成功，那么 *pattern* 匹配的字符串必须开始于 *string* 的最左端。也就是说，re.match()只会在字符串的开始位置尝试匹配，re.search()则没有这个限制。

```
re.search(r"\d", "ab1") != None      # => True
re.match(r"\d", "ab1") != None      # => False
re.search(r"\d", "1ab") != None     # => True
re.match(r"\d", "1ab") != None     # => True
```

`re.match()` 的名字容易让人认为这是专门用来进行数据验证的方法。其实单纯调用它并不能准确验证，还必须在正则表达式末尾加上 `\Z`，不过更好的做法是在正则表达式首尾加上 `\A` 和 `\Z`，这样更符合惯例。

```
re.search(r"\d{6}", "123456") != None      # => True
#错误的验证
re.search(r"\d{6}", "123456ab") != None    # => True
#尾部加上\Z
re.search(r"\d{6}\Z", "123456ab") != None  # => False
#推荐的做法
re.search(r"\A\d{6}\Z", "123456ab") != None # => True
```

### 14.3.6 re.findall(pattern, string[, flags])

这个函数会一次性找出正则表达式 *pattern* 在 *string* 中的所有匹配，返回一个 `list`，最简单的用法就是：

```
re.findall(r"\d{4}-\d{2}-\d{2}", "2010-12-20 2011-02-14")
['2010-12-20', '2011-02-14']
```

如果 *pattern* 中存在捕获分组，则返回 `list` 的元素并非整个表达式所匹配的文本，而是各个捕获分组所匹配文本构成的 `tuple`。

```
re.findall(r"(\d{4})-(\d{2})-(\d{2})", "2010-12-20 2011-02-14")
[('2010', '12', '20'), ('2011', '02', '14')]
```

请注意，如果 *pattern* 中存在捕获分组，整个表达式匹配的文本不会包含在返回文本中，而且返回 `tuple` 中的元素偏移值与分组编号并不对应（编号为 1 的分组匹配的文本是 `tuple` 中的元素 0，编号为 2 的分组匹配的文本是 `tuple` 中的元素 1，依此类推）。如果需要得到整个表达式匹配的文本，则可以显式给整个表达式增加括号。

```
re.findall(r"((\d{4})-(\d{2})-(\d{2}))", "2010-12-20 2011-02-14")
[('2010-12-20', '2010', '12', '20'), ('2011-02-14', '2011', '02', '14')]
```

因为它返回的是 `list` 而不是 `dict`，所以即便使用了命名分组，`re.findall()` 也只会将它们当做数字编号分组来处理，而不会保存任何与命名有关的信息。

```
regex = r"((?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2}))"
re.findall(regex, "2010-12-20 2011-02-14")
[('2010-12-20', '2010', '12', '20'), ('2011-02-14', '2011', '02', '14')]
```

### 14.3.7 re.finditer(pattern, string[, flags])

如果文本内容很多，使用 `re.findall()` 一次性找出所有匹配的成本可能很高，使用

`re.finditer()`可以返回一个迭代器 (iterator)，它可以按从左到右的顺序逐个遍历每次匹配的 `MatchObject` 对象，依次进行处理或检验。

```
#看起来是先一次性找到所有匹配结果，再遍历，其实每次的匹配结果是临时得到的
regex = r"(\d{4})-(\d{2})-(\d{2})"
for iter in re.finditer(regex, "2010-12-20 2011-02-14"):
    print "%s, %s, %s, %s" % (iter.group(), iter.group(1), iter.group(2), iter.group(3))

2010-12-20, 2010, 12, 20
2011-02-14, 2011, 02, 14

#使用命名分组
regex = r"((?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2}))"
for iter in re.finditer(regex, "2010-12-20 2011-02-14"):
    print "%s, %s, %s, %s" % (iter.group(), iter.group(1), iter.group(2), iter.group(3))

2010-12-20, 2010, 12, 20
2011-02-14, 2011, 02, 14
```

### 14.3.8 re.split(pattern, string[, maxsplit=0, flags=0])

这个函数用正则表达式 *pattern* 能够匹配的文本切分字符串 *string*。

```
re.split(r"-", "2010-12-20")
['2010', '12', '20']
```

如果正则表达式在字符串的开头或结尾位置能找到匹配，那么结果的开头或结尾会出现空字符串。

```
re.split(r"-", "-2010-12-20-")
['', '2010', '12', '20', '']
```

也可以明确设定 *maxsplit*。请注意，在其他语言中，这个参数一般表示返回数组的长度，但是在 Python 中，它表示切分的次数。也就是说，返回数组一般会包含 *maxsplit*+1 个元素。如果 *maxsplit* 设定为负数，则表示“不做任何切分”；若设定为 0，则等于没有指定；若设定的值为正数，且小于实际可切分次数，则按设定的值进行切分，若大于实际可切分次数，则以实际可切分次数为准。

```
re.split(r"-", "2010-12-20", -1)
['2010-12-20']

re.split(r"-", "2010-12-20", 0)
['2010', '12', '20']

re.split(r"-", "2010-12-20", 1)
['2010', '12-20']
```

### 14.3.9 re.sub(pattern, repl, string[, count, flags])

这个函数用来进行正则表达式替换，它将字符串 *string* 中正则表达式 *pattern* 能够匹配的文本替换为 *repl* 指定的文本。

```
regex = r"\d{4}-\d{2}-\d{2}"
replacement = "Date"
re.sub(regex, replacement, "2010-12-20 2011-02-14")
Date Date
```

如果要在 *replacement* 字符串中引用之前某个捕获分组匹配的文本，则应当使用 `\num`，不过因为 *replacement* 是字符串，如果使用字符串文字，则应当写成 `\\num`。为了清晰直接，推荐使用原生字符串，也就是在字符串之前添加 `r`。

```
regex = r"(\d{4})-(\d{2})-(\d{2})"
replacement = r"\2/\3/\1"
print re.sub(regex, replacement, "2010-12-20")
12/20/2010
```

Python 也提供了更清晰的方式在 *repl* 中引用分组，写作 `\g<num>`。

```
regex = r"(\d{4})-(\d{2})-(\d{2})"
replacement = r"\g<2>/\g<3>/\g<1>"
print re.sub(regex, replacement, "2010-12-20")
12/20/2010
```

这种记法也可以用来引用命名分组，写作 `\g<name>`。

```
regex = r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})"
replacement = r"\g<month>/\g<day>/\g<year>"
print re.sub(regex, replacement, "2010-12-20")
12/20/2010
```

*repl* 也可以是一个函数，它接收一个 `MatchObject` 对象，返回一个 `String` 对象。这样替换功能就强大了许多，比如把所有单词统一为首字母大写格式。

```
def capitalize(match) :
    #正则表达式已经用两个捕获分组分开第一个字母和之后的字母
    return match.group(1).upper() + match.group(2).lower()

print re.sub(r"(?i)\b([a-z])([a-z]+)\b", capitalize, "one TWO three ")
One TWO THREE
```

也可以设定 *count* 参数，它指定替换操作最多能发生的次数。

```
print re.sub(r"\bt[a-z]+\b", toUpper, "one two three", 1)
one TWO three
```

## 14.4 常用操作示例

### 14.4.1 验证

#### 简单验证

```
#应当在正则表达式首尾加上\A 和\Z
re.search(r"\A\d{4}-\d{2}-\d{2}\Z", "2010-12-20") != None # => True
```

#### 复用 RegexObject 验证

```
dateRegex = re.compile(r"\A\d{4}-\d{2}-\d{2}\Z")
re.search(dateRegex, "2010-12-20") != None # => True
```

### 14.4.2 提取

#### 简单提取

```
regex = r"\d{4}-\d{2}-\d{2}"
string = "2010-12-20 2011-02-14"
for dateStr in re.findall(r"\d{4}-\d{2}-\d{2}", "2010-12-20 2011-02-14") :
    print dateStr
2010-12-20
2011-02-14
```

#出现捕获分组之后，返回的结果不会包含默认的捕获分组 0，所以显式加上分组 0

```
regex = r"((\d{4})-(\d{2})-(\d{2}))"
string = "2010-12-20 2011-02-14"
for matchTuple in re.findall(regex, string) :
    print "date:", matchTuple[0],
    print "year:", matchTuple[1],
    print "month:", matchTuple[2],
    print "day:", matchTuple[3]
date: 2010-12-20 year: 2010 month: 12 day: 20
date: 2011-02-14 year: 2011 month: 02 day: 14
```

#### 复用 RegexObject 对象

```
regex = re.compile(r"((\d{4})-(\d{2})-(\d{2}))")
string = "2010-12-20 2011-02-14"
for matchTuple in regex.findall(string) :
    print "date:", matchTuple[0],
    print "year:", matchTuple[1],
    print "month:", matchTuple[2],
    print "day:", matchTuple[3]
```

```
date: 2010-12-20 year: 2010 month: 12 day: 20
date: 2011-02-14 year: 2011 month: 02 day: 14
```

#### 使用 finditer 迭代查找

```
regex = re.compile(r"(\d{4})-(\d{2})-(\d{2})")
string = "2010-12-20 2011-02-14"
for iter in regex.finditer(string) :
    print "date:", iter.group(0),
    print "year:", iter.group(1),
    print "month:", iter.group(2),
    print "day:", iter.group(3)
date: 2010-12-20 year: 2010 month: 12 day: 20
date: 2011-02-14 year: 2011 month: 02 day: 14
```

#### 使用命名分组

```
regex = re.compile(r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})")
string = "2010-12-20 2011-02-14"
for iter in regex.finditer(string) :
    print "date:", iter.group(),
    print "year:", iter.group("year"),
    print "month:", iter.group("month"),
    print "day:", iter.group("day")
date: 2010-12-20 year: 2010 month: 12 day: 20
date: 2011-02-14 year: 2011 month: 02 day: 14
```

### 14.4.3 替换

#### 简单替换

```
regex = r"(\d{4})-(\d{2})-(\d{2})"
replacement = r"\g<2>/\g<3>/\g<1>"
print re.sub(regex, replacement, "2010-12-20")
12/20/2010
```

#### 命名分组替换

```
regex = r"(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})"
replacement = r"\g<month>/\g<day>/\g<year>"
print re.sub(regex, replacement, "2010-12-20")
12/20/2010
```

#### 在替换中使用整个表达式匹配的文本

```
regex = r"\d+\.\d{0,2}"
replacement = r"$\g<0>"
```

```
print re.sub(regex, replacement, "the price is 12.99")
the price is $12.99
```

#### 14.4.4 切分

```
regex = r"\s+"
for word in re.split(regex, "one\ntwo\n three") :
    print word
one
two
three
```

## 第 15 章 Ruby

Ruby 内建了对正则表达式的支持，提供了相对完备的正则表达式功能。但 Ruby 1.8 的正则表达式也存在不少限制，最突出的是不支持逆序环视。Ruby 1.9 采用全新的 Oniguruma（鬼车）正则引擎，补充了包括逆序环视在内的许多功能，最重要的是全面支持 Unicode（但是也造成了一些奇怪的现象，下面你会看到）。不过考虑到目前 Ruby 1.9 并没有完全取代 Ruby 1.8，所以本章中我们会讲解 1.8 和 1.9 的差异，章末的表格详细列出了 Ruby 1.9 中正则表达式的功能变化，有兴趣的读者可以参考<sup>1</sup>。

### 15.1 预备知识

Ruby 的正则表达式对象是 `Regexp`<sup>2</sup>。最常见的生成方法类似 `/regex/`，即在首尾两个斜线/之间，用正则文字给出正则表达式，比如 `/\d+/`。

为继续讲解方便，先介绍进行正则表达式匹配的方法：Ruby 中的正则匹配非常简单，直接使用操作符  `=~`  连接正则表达式与字符串即可。如果可以匹配，则返回第一次匹配发生的起始位置，否则返回 `nil`。请注意，正则表达式和字符串的左右顺序并不影响结果。

```
"0" =~ /\d/ # =>0
"a" =~ /\d/ # =>nil
/\d/ =~ "0" # =>0
/\d/ =~ "a" # =>nil
```

Ruby 中有两种字符串：双引号字符串和单引号字符串。两者的主要区别在于：双引号字符串需要进行字符串转义处理，而单引号字符串不需要进行字符串转义处理<sup>3</sup>。所以如果要作为正则表达式的字符串，最好使用单引号字符串，这样就可以直接使用 `\d`，而不要写成 `\\d`；作为普通文本的字符串，则最好使用双引号字符串。

---

<sup>1</sup> 本书中所说的 Ruby 1.9 指的是 Ruby 1.9.1 及以上版本，而不包括 Ruby 1.9.0，因为这个版本用的人极少，而且经过作者测试，Ruby 1.9.0 的正则表达式在不同操作系统上的表现差异非常大。如果你使用的是 Ubuntu 之类的操作系统，通过 `apt-get` 安装的 Ruby 1.9，很可能是 Ruby 1.9.0，请注意这一点。

<sup>2</sup> 如果用过 Ruby 中的正则表达式，可能会发现 Ruby 和 JavaScript 中的正则表达式对象名字很像，唯一的区别在于 Ruby 的 `Regexp` 中 `e` 是小写，而 JavaScript 的 `RegExp` 中 `E` 是大写。

<sup>3</sup> 关于字符串和转义的详细解释，请参考本书第 6 章。



## 15.2 正则功能详解

### 15.2.1 列表

表 15-1 中列出了 Ruby 的正则功能。

表 15-1 Ruby 的正则功能

功能	记法	说明
字符组 ④1	<code>[...]</code> <code>[^...]</code>	支持 Unicode
POSIX 字符组 ④15	<code>[:digit:]</code>	1.8 只能匹配 ASCII 字符, 1.9 可以匹配 Unicode 字符
Unicode 属性 ④121	<code>\p{...}</code>	1.9 支持
字符组简记法 ④115	<code>\d</code> <code>\D</code> <code>\w</code> <code>\W</code> <code>\s</code> <code>\S</code>	采用 ASCII 匹配规则
行起始位置 ④61	<code>^</code> <code>\A</code>	默认采用多行模式
行结束位置 ④61	<code>\$</code> <code>\z</code> <code>\Z</code>	默认采用多行模式
单词边界 ④117	<code>\b</code> <code>\B</code>	1.8 采用 ASCII 匹配规则, 1.9 采用 Unicode 匹配规则
顺序环视 ④68	<code>(?=...)</code> <code>(?!...)</code>	支持
逆序环视 ④68	<code>(?&lt;=...)</code> <code>(?&lt;!...)</code>	1.9 支持
匹配模式 ④81	<code>i</code> <code>m</code> <code>o</code> <code>x</code> <code>u</code>	支持
模式作用范围 ④89	<code>(?modifier)</code> <code>(?-modifier)</code>	支持
纯文本模式 ④99	<code>\Q...\E</code>	不支持
捕获分组及引用④44	<code>(...)</code> <code>\num</code> <code>\$num</code>	支持
命名分组 ④53	<code>(?&lt;name&gt;...)</code> <code>\k&lt;name&gt;</code>	1.9 支持
非捕获分组 ④54	<code>(?:...)</code>	支持
多选结构 ④39	<code>(... ...)</code>	支持
匹配优先量词 ④18	<code>?</code> <code>*</code> <code>+</code>	支持
忽略优先量词 ④26	<code>??</code> <code>?*</code> <code>?+</code> <code>{n,m}?</code>	1.9 支持

### 15.2.2 字符组

Ruby 中的正则引擎可以正确识别 ASCII 编码, 但无法正确识别 GBK 编码, 因此无法避免 GBK 编码的“错误匹配”问题 (④111)。

"罚" =~ /[收发]/ # => 0

解决的办法是使用 Unicode 编码环境, 指定程序使用 Unicode 编码的办法是在源文件的顶部

写上这样一行：

```
# encoding:UTF-8
```

Ruby 1.9 之后版本的正则表达式才提供了相对完备的 Unicode 支持，如果你使用的是 Ruby 1.8，正则表达式中又有非 ASCII 字符，则最好在表达式末尾写上 `u` 显式指定 Unicode 模式(Ruby1.9 中正则表达式默认使用 Unicode 编码，所以不用显式指定)。

```
"罚" =~ /[收发]/u # => nil
"发" =~ /[收发]/u # => 0
```

另外要注意的是，在 Ruby 1.8 中，即便指定程序使用 Unicode 编码，匹配位置仍然是按照字节计算的，Ruby 1.9 中则是按照字符计算的。

#### Ruby 1.8

#UTF-8 编码每个汉字占 3 个字节，所以匹配的“发”从 3 开始

```
"我发" =~ /[收发]/u # => 3
```

#### Ruby 1.9

```
"我发" =~ /[收发]/ # => 1
```

在 Ruby 中可以用字符组 `[\u{4e00}-\u{9fff}]` 匹配任意中文字符，需要注意的是，这种写法只有 Ruby 1.9 以上版本才支持，并且使用时必须显式指定 Unicode 模式。

```
"我" =~ /[\u{4e00}-\u{9fff}]/u # => 0
"你" =~ /[\u{4e00}-\u{9fff}]/u # => 0
"a" =~ /[\u{4e00}-\u{9fff}]/u # => nil
```

Ruby 支持 POSIX 字符组。在 Ruby 1.8 中，POSIX 字符组只能匹配 ASCII 字符(指定 Unicode 模式也不行)，而且不可使用 `[:ascii:]` 和 `[:word:]`；而在 Ruby 1.9 中，POSIX 字符组可以匹配 Unicode 字符，同时支持 `[:ascii:]` 和 `[:word:]`。

#### Ruby 1.8

```
"0" =~ /[[:digit:]]/ # => 0
"a" =~ /[[:alpha:]]/ # => 0
```

#### Ruby 1.9

```
"0" =~ /[[:ASCII:]]/ # => 0
"a" =~ /[[:word:]]/ # => 0
"我" =~ /[[:word:]]/ # => 0
#全角数字
"0" =~ /[[:digit:]]/ # => 0
```

### 15.2.3 Unicode 属性

Ruby 1.9 支持使用 Unicode Script，只是使用 Unicode Script 时，必须显式指定 `u` 模式，否则会报错。所以，如果要匹配中文字符，则可以使用 `\p{Han}`。

```
"a" =~ /\p{Han}/u # => nil
```

```
"我" =~ /\p{Han}/u      # => 0
"我" =~ /\p{Han}/      # 报错
```

## 15.2.4 字符组简记法

Ruby 的字符组简记法都采用 ASCII 编码匹配规则。如果指定了 Unicode 模式，则会换用 Unicode 匹配规则<sup>1</sup>。

```
#半角字符
"1"  =~ /\d/      # => 0
"a"  =~ /\w/      # => 0
" "  =~ /\s/      # => 0
#全角字符及中文字符
"1 " =~ /\d/      # => nil
"我" =~ /\w/      # => nil
" "  =~ /\s/      # => nil
```

## 15.2.5 单词边界

在 Ruby 1.8 中，如果没有指定 Unicode 模式，`\w` 采用 ASCII 匹配规则，单词边界 `\b` 也采用这种匹配规则（☞115）。

```
#英文字符-半角标点
"a," =~ /a\b/      # => 0
#英文字符-全角标点
"a, " =~ /a\b/      # => 0
#英文字符-结束
"a"  =~ /a\b/      # => 0
#中文字符-半角标点
"我," =~ /我\b/     # => nil
#中文字符-全角标点
"我, " =~ /我\b/     # => nil
#中文字符-空字符串
"我"  =~ /我\b/     # => nil
#全角标点-结束
", "  =~ /\b/       # => nil
#中文字符-英文字符
"我a" =~ /我\b/     # => 0
```

如果指定了 Unicode 模式，则 `\b` 也采用 Unicode 匹配规则。

```
#英文字符-半角标点
"a," =~ /a\b/u      # => 0
```

---

<sup>1</sup> 如果你使用的是 Ruby 1.9.0，测试结果可能与这里的完全不同。如本章开头所述，Ruby 1.9.0 是一个非常奇怪的版本，不推荐使用。

```

#英文字符-全角标点
"a, " =~ /a\b/u # => 0
#英文字符-结束
"a" =~ /a\b/u # => 0
#中文字符-半角标点
"我," =~ /我\b/u # => 0
#中文字符-全角标点
"我, " =~ /我\b/u # => 0
#中文字符-空字符串
"我" =~ /我\b/u # => 0
#全角标点-结束
", " =~ /\b/u # => 0
#中文字符-英文字符
"我 a" =~ /我\b/u # => 0

```

Ruby 1.9 中的 `\w` 虽然采用 ASCII 匹配规则，但 `\b` 的表现与 Ruby 1.8 中指定 Unicode 方式一样，这一点值得特别注意。

### 15.2.6 行起始/结束位置

Ruby 默认采用多行模式，所以 `^` 不但可以匹配整个字符串的开始位置，还可以匹配字符串内部的文本行的起始位置；而 `\A` 无论在什么情况下都匹配整个字符串的起始位置。

```

"1\n2\n" =~ /^1/ # => 0
"1\n2\n" =~ /^2/ # => nil

"1\n2\n" =~ /\A1/ # => 2
"1\n2\n" =~ /\A2/ # => nil

```

同样，因为默认采用多行模式，所以 `$` 可以匹配文本内部行的结束位置，`\Z`、`\z` 匹配的是整个字符串的结束位置（如果结束位置有换行符，则 `\Z` 匹配这个换行符之前的位置）。

```

"1\n2\n" =~ /1$/ # => 0
"1\n2\n" =~ /2$/ # => 2
"1\n2" =~ /2$/ # => 2

"1\n2\n" =~ /1\Z/ # => nil
"1\n2\n" =~ /2\Z/ # => 2
"1\n2" =~ /2\Z/ # => 2

"1\n2\n" =~ /2\z/ # => nil
"1\n2" =~ /2\z/ # => 2

```

重复一次，Ruby 的 `m` 模式不等于其他语言中的多行模式，而是等于通常的单行模式，它只影响点号 `.` 的匹配。

15.2.7 环视

Ruby 中的肯定顺序环视 `(?=...)` 和否定顺序环视 `(?!...)` 内可以使用任意形式的子表达式。

```
"cab" =~ /c(?=ab+)/ # => 0
"ccd" =~ /c(?(ab|cd))/ # => 0
"cabbbb" =~ /c(?(ab+|cd))/ # => 0
```

对逆序环视的支持则不那么好：Ruby 1.8 完全不支持逆序环视；Ruby 1.9 支持逆序环视，但是逆序环视中的表达式匹配的文本长度必须是固定的，完全不能使用 `+`、`*`、`?` 之类的量词：如果使用了多选结构，只能出现竖线，比如 `regex1|regex2`，而不能出现括号，比如 `(regex1|regex2)`。

```
"ab" =~ /ab(?!ab)/ # => 0
"cd" =~ /cd(?!ab|cde)/ # => 0

#ab?能匹配的文本长度不固定，报错
"ab" =~ /ab(?!ab?)/

#多选结构(ab|cd)包含括号，报错
"ab" =~ /ab(?!=(ab|cd))/
```

15.2.8 匹配模式

表 15-2 中列出了 Ruby 的正则表达式支持的匹配模式。

表 15-2 Ruby 的正则表达式支持的匹配模式

常量	修饰符	解释
Regexp::IGNORECASE	i	不区分 ASCII 字符的大小写
Regexp::EXTENDED	x	容许正则表达式中出现注释，表达式中的空白字符，以及# 开始到行末的文本，都视为注释
Regexp::MULTILINE	m	等于通常所说的“单行模式”，容许点号匹配任何字符，包括换行符
	u	对字符组简记法采取 Unicode 匹配规则
	o	对包含插值 ( <code>#{variable}</code> ) 的正则表达式，只在第一次遇到时求值，以后使用时不变；否则每次遇到时都求值

注：修饰符 `u` 其实是用来指定正则表达式编码的，可选的值一共有 4 个：`n`，表示不指定；`e`，表示 EUC (Extended UNIX Code，一种主要用于中日韩文字的多字节编码)；`u`，表示 Unicode；`s`，表示 SJIS (Shift JIS，一种主要用于日文的多字节编码)。通常只会用到 `u`。

匹配模式可以在表达式中以 `(?modifier)` 指定，或者在结束分隔符之后写上模式对应的修

饰符，也可以在生成 `Regexp` 对象时，通过预定义常量作为参数指定。

```
"A" =~ /a/ # => nil
"A" =~ /(?)a/ # => 0
"A" =~ /a/i # => 0
"A" =~ Regexp.new('a', Regexp::IGNORECASE) # => 0
```

Ruby 支持用 `(?-modifier)` 停用某个模式。

```
"aBb" =~ /a(?!b)(?-i)b/ # => 0
"aBB" =~ /a(?!b)(?-i)b/ # => nil
```

### 15.2.9 捕获分组的引用

在 Ruby 中，如果要在正则表达式内部引用捕获分组，则应当使用 `\num` 记法，其中 `num` 为对应捕获分组的编号。

```
"ab" =~ /([a-z])\1/ # => nil
"aa" =~ /([a-z])\1/ # => 0
```

如果要在替换时引用捕获分组，则应当使用 `$num` 记法。

```
"2010-12-20".sub(/(\d{4})-(\d{2})-(\d{2})/, '\2/\3/\1')
12/20/2010
```

如果使用了命名分组，在表达式内部则应当使用 `\k<name>` 来引用。

```
"ab" =~ /(?!<char>[a-z])\k<char>/ # => nil
"aa" =~ /(?!<char>[a-z])\k<char>/ # => 0
```

在 *replacement* 字符串中同样使用 `\k<name>` 来引用。

```
"2010-12-20".sub(/(?<year>\d{4})-(?<month>\d{2})-(?<day>\d{2})/,
'\k<month>/\k<day>/\k<year>')
12/20/2010
```

## 15.3 正则 API 简介

### 15.3.1 Regexp

Ruby 中的正则表达式对象是 `Regexp`，一般使用 `/regex/` 来生成。如果习惯更为“正统”的模式，也可以用 `Regexp.new()` 来生成，参数是表达式对应的字符串。

```
Regexp.new('[a-z]') =~ "a" # => 0
```

之前讲过，Ruby 中有两种字符串，其中双引号字符串需要进行字符串转义处理，而单引号

字符串不需要进行字符串转义处理。所以，作为正则表达式的字符串，最好使用单引号字符串。

如果要指定匹配模式，可以添加第 2 个参数（虽然实际参数是一个整数，但是 Ruby 其实有 3 个预定义常量，即 `Regexp::IGNORECASE`、`Regexp::EXTENDED`、`Regexp::MULTILINE`）。

```
Regexp.new('[a-z]', Regexp::IGNORECASE) =~ "a"      # => 0
Regexp.new('[a-z] #comment', Regexp::EXTENDED) =~ "a" # => 0
Regexp.new('.', Regexp::MULTILINE) =~ "\n"          # => 0
```

也可以混用多个模式，用|连接预定义常量。

```
Regexp.new('[a-z]#comment', Regexp::IGNORECASE|Regexp::EXTENDED) =~ 'A' # => 0
```

`Regexp.new()` 的第 1 个参数也可以不用字符串，而是用类似 `/regex/` 的形式。

```
Regexp.new(/[a-z]/) =~ "a" # => 0
```

如果使用这种形式，那么有三种形式指定匹配模式，下面各举一个例子，其结果相同。

```
Regexp.new(/(?i)[a-z]/) =~ "a" # => 0
Regexp.new(/[a-z]/i) =~ "a"    # => 0
Regexp.new(/[a-z]/, Regexp::IGNORECASE) =~ "a" # => 0
```

另外还有 `Regexp.compile()` 方法，它与 `Regexp.new()` 完全一致。

值得注意的还有 `Regexp` 对象的编码。在 Ruby 1.9 中，`Regexp` 对象自身是具有编码的，在默认情况下，如果正则表达式中出现的都是 ASCII 字符，且程序采用了兼容 ASCII 的编码，则正则表达式默认采用 ASCII 编码；否则，采用与程序相同的编码。调用 `Regexp.encoding?` 可以得到编码信息。

```
/abc/.encoding      #=> #<Encoding:ASCII>

#如果设定了源代码使用 UTF-8 编码
/你我他/.encoding   #=> #<Encoding:UTF-8>

#也可以显式指定 Regexp 的编码
/abc/u.encoding     #=> #<Encoding:UTF-8>
```

要正确使用具有编码属性的正则表达式，只能是两种情况：`Regexp` 与字符串使用相同的编码；或者 `Regexp` 使用 ASCII 编码，而字符串使用兼容 ASCII 的编码，否则会报错。

```
Regexp.new('a'.encode('iso-8859-1')) =~ "你我他"
#=> Encoding::CompatibilityError: incompatible encoding regexp match

(ISO-8859-1 regexp with UTF-8 string)
```

如果不确定使用的 `Regexp` 对象的编码是否兼容 ASCII，可以使用 `fix_encoding?` 来判断。不兼容的情况下会返回 `true`；如果一个 `Regexp` 对象的编码兼容 ASCII，它可以尝试匹配任何使用了 ASCII 编码的字符串，而不会报错。

```
Regexp.new('你我'.encode('gbk')).fixed_encoding? # => True
Regexp.new('你我'.encode('utf-8')).fixed_encoding? # => True
```

### 15.3.2 Regexp.match(text)

之前看到的匹配都是依靠操作符`=~`进行的，许多读者可能不熟悉，所以这里先详细介绍它。`=~`用来执行正则表达式和字符串的匹配，返回匹配成功的起始位置；如果匹配不成功，则返回`nil`。因此，可以用这个办法来进行数据验证（验证数据时，别忘了在正则表达式两端加上`\A`和`\Z`）。

```
("2010-12-20" =~ /\A\d{4}-\d{2}-\d{2}\Z/) != nil # => True
("2010-12" =~ /\A\d{4}-\d{2}-\d{2}\Z/) != nil # => False
("a2010-12-20" =~ /\A\d{4}-\d{2}-\d{2}\Z/) != nil # => False
```

如果匹配成功，会生成一个`MatchData`对象，通过`Regexp.last_match`可以获得它，其中包含了匹配的详细信息。

```
"Date: 2010-12-20, " =~ /(\d{4})-(\d{2})-(\d{2})/
#获得 MatchData 对象
data = Regexp.last_match
# => #<MatchData "2010-12-20" 1:"2010" 2:"12" 3:"20">
#正则表达式中捕获分组数目+1，因为整个表达式匹配的文本作为编号 0
data.size() # => 4
#第 0 号分组，也就是整个表达式匹配的文本
data[0] # => "2010-12-20"
#第 1、2、3 号分组，各捕获分组匹配的文本
data[1] # => "2010"
data[2] # => "12"
data[3] # => "20"
#编号最大的捕获分组匹配的文本
data[-1] # => "20"
#各捕获分组匹配的文本，编号从 1 开始
data.captures # => ["2010", "12", "20"]
#编号为 1 的捕获分组的开始位置、结束位置、范围
data.begin(1) # => 0
data.end(1) # => 4
data.offset(1) #[0, 4]
#原字符串中，表达式所匹配文本左侧的文本
data.pre_match # => "Date: "
#原字符串中，表达式所匹配文本右侧的文本
data.post_match # => ", "
```

也可以不用显式通过`MatchData`对象来获得这些信息，Ruby 提供了若干特殊变量，使用它们同样也可以达到目的。表 15-3 列出了常用的特殊变量。



表 15-3 Ruby 提供的特殊变量

特殊变量	等价
<code>\$~</code>	<code>Regexp.last_match</code>
<code>\$&amp;</code>	<code>Regexp.last_match[0]</code>
<code>\$`</code>	<code>Regexp.last_match.pre_match</code>
<code>\$'</code>	<code>Regexp.last_match.post_match</code>
<code>\$1</code>	<code>Regexp.last_match[1]</code>
<code>\$2</code>	<code>Regexp.last_match[2]</code>
<code>...</code>	<code>...</code>
<code>\$+</code>	<code>Regexp.last_match[-1]</code>

`Regexp.match(text)`方法与 `=~` 差不多，主要的区别在于，它返回的并不是匹配的起始位置，而是一个 `MatchData` 对象，如果匹配不成功，则返回 `nil`。所以，上面的代码可以这样获得 `MatchData`。

```
data = /\A(\d{4})-(\d{2})-(\d{2})\z/.match("2010-12-20")
```

`Regexp.match(string)`与`=~`的另一不同是，它可以通过可选参数设定偏移值，设定尝试匹配的起始位置，这样就可以在字符串中逐次找到匹配。在下面的例子中，每次匹配成功之后，记录下匹配的结束位置，作为下次开始尝试的起点，就实现了逐次匹配。

```
regex = /\d{4}-\d{2}-\d{2}/
text = "2010-12-20 2011-02-14"
data = regex.match(text)
while data != nil
  puts data[0]
  data = regex.match(text, data.end(0))
end
2010-12-20
2011-02-14
```

15.3.3 Regexp.quote(text)和 Regexp.escape(text)

这两个方法是等价的，它们消除 `text` 中可能出现的任何元字符的特殊含义。

```
Regexp.new(Regexp.quote("[a-z]")) =~ "[a-z]" # => 0
Regexp.new(Regexp.quote("[a-z]")) =~ "a"      # => nil
Regexp.new(Regexp.quote("[a-z]")) =~ "z"      # => nil

Regexp.new(Regexp.escape("[a-z]")) =~ "[a-z]" # => 0
Regexp.new(Regexp.escape("[a-z]")) =~ "a"      # => nil
Regexp.new(Regexp.escape("[a-z]")) =~ "z"      # => nil
```

### 15.3.4 String.index(Regexp)

这个方法默认等价于`==`操作，如果匹配成功，则返回匹配开始的位置，否则返回 `nil`。

如果不指定第 2 个参数，它等价于 `Regex.match()` 或者 `==`。

```
"0".index(/\d/)    # => 0
"a0".index(/\d/)   # => 1
"a".index(/\d/)    # => nil
```

它也可以指定第 2 个参数，设定开始尝试匹配的偏移值。此时可以遍历文本，获得逐次匹配的结果。

```
regex = /\d{4}-\d{2}-\d{2}/
text = "2010-12-20 2011-02-14"
offset = text.index(regex)
while offset != nil
  puts Regexp.last_match[0]
  offset = text.index(regex, Regexp.last_match.end(0))
end

2010-12-20
2011-02-14
```

### 15.3.5 String.scan(Regexp)

`==`操作符和 `Regexp.match(string)` 方法只能进行单次匹配，如果希望一次性找到某个正则表达式在字符串中的所有匹配，就应该使用 `String.scan(Regexp)`，它返回各次匹配所对应文本组成的数组。

```
puts "2010-12-20 2011-02-14".scan(/\d{4}-\d{2}-\d{2}/)
["2010-12-20", "2011-02-14"]
```

如果正则表达式包含捕获分组，则返回数组的元素本身也是数组。

```
puts "2010-12-20 2011-02-14".scan(/((\d{4})-(\d{2})-(\d{2}))/)
[["2010", "12", "20"], ["2011", "02", "14"]]
```

请注意，这时候并不会返回整个表达式匹配的文本（也就是编号为 0 的捕获分组捕获的文本），如果需要得到整个表达式匹配的文本，可以给整个表达式首尾加上括号。

```
puts "2010-12-20 2011-02-14".scan(/((\d{4})-(\d{2})-(\d{2})))/)
[["2010-12-20", "2010", "12", "20"], ["2011-02-14", "2011", "02", "14"]]
```

也可以用 `String.scan()` 进行更复杂的字符串处理，方法是在括号后面加上 `{|match, ...|block}`，其中 `match,...` 表示各捕获分组匹配的文本，而 `block` 则是对这些文本进行处理的程序代

码，比如下面这样。

```
#匹配同时将单词转大写输出
"one two three".scan(/\b\w+\b/) {|x| puts x.upcase}

ONE

TWO

THREE

one two three
```

### 15.3.6 String.slice(Regexp)

这个方法返回 *Regexp* 在 *String* 中能匹配的文本。

```
"a0".slice(/\d/) # => "0"
```

也可以设定第 2 个参数，指定对应捕获分组的编号，则返回对应的捕获分组匹配的文本。

```
string = "Date: 2010-12-20"
regex = /(\d{4})-(\d{2})-(\d{2})/
string.slice(regex, 1) # => "2010"
string.slice(regex, 2) # => "12"
string.slice(regex, 3) # => "20"
```

如果正则表达式可以在文本中找到多次匹配，那么 **String.slice(Regexp)** 只返回第一次匹配的结果。

```
string = "Date: 2010-12-20"
regex = /\d{4}-\d{2}-\d{2}/
string.slice(regex) # => 2010-12-20
```

如果将 **String.slice(Regexp)** 写成 **String.slice!(Regexp)**，则 *String* 自身也受到影响：在返回匹配文本的同时，还会从原字符串中删除匹配的结果。

```
text = "2010-12-20 2011-02-14"
text.slice!(/\d{4}-\d{2}-\d{2}/)
puts text
2011-02-14
```

### 15.3.7 String.split(Regexp)

这个方法用正则表达式来切分字符串，它返回一个数组，其中的元素是切分之后的文本。

```
"2010-12-20".split(/-/)
["2010", "12", "20"]
```

如果切分的末尾留下空字符串，在默认情况下会忽略它，而字符串开头的空字符串不会忽略。

```
"-2010-12-20-".split(/-/)  
["", "2010", "12", "20"]
```

这个方法也可以指定第 2 个参数，它是一个整数，假设为  $n$ ，按照  $n$  取值的不同，有几种情况，以下分别举例说明。

```
# n<0, 切分尽量多次，同时保留最后的空字符串；  
"-2010-12-20-".split(/-/, -1)  
["", "2010", "12", "20", ""]  
  
# n=0, 等于未设定此参数，切分尽量多次，不保留最后的空字符串；  
"-2010-12-20-".split(/-/, 0)  
["", "2010", "12", "20", ""]  
  
# n<=limit, 等于未设定此参数，切分尽量多次，不保留最后的空字符串；  
"-2010-12-20-".split(/-/, 2)  
["", "2010-12-20-"]  
  
# n > limit, 进行 n 次切分，末尾的空字符串会保留。  
"-2010-12-20-".split(/-/, 10)  
["", "2010", "12", "20", ""]
```

### 15.3.8 String.sub(Regexp, Str)

`String.sub(Regexp, Str)` 用来执行正则表达式替换，*Regexp* 是匹配想要替换文本的表达式，而 *Str* 是希望更新的字符串，调用之后只会对 *Regexp* 第一次能匹配的文本进行替换。

```
puts "2010-12-20 2011-02-14".sub(/\d{4}-\d{2}-\d{2}/, 'date')  
date 2011-02-14
```

如果第 1 个参数不是 *Regexp* 而是普通字符串，则进行普通字符串替换。

```
puts ".* 2011-02-14".sub('.*', 'date')  
date 2011-02-14
```

也可以在 *String* 参数中，用  $\backslash num$  引用 *Regexp* 对应捕获分组所匹配的文本。

```
puts "2010-12-20".sub(/(\d{4})-(\d{2})-(\d{2})/, '\2/\3/\1')  
12/20/2010
```

如果需要进行更灵活的替换，则可以使用 `String.sub(Regexp){| match | block}` 方式，其中的 *match* 是 *Regexp* 所匹配的文本，即一个 *String* 对象，我们可以写一段代码作为 *block*，进行更复杂的处理，比如将单词转换为大写。

```
#注意只会转换处理第一次匹配的文本  
puts "one two three".sub(/[a-z]+/){| s | s.upcase}  
ONE two three
```

如果需要在 *block* 中获得其他匹配的内容，可以使用之前讲过的 `$1`、`$2`、`$`` 等特殊变量。

```
regex = /(\d{4})-(\d{2})-(\d{2})/
puts "2010-12-20".sub(regex){|s| $2 + "/" + $3 + "/" + $1}
12/20/2010
```

`String.sub(Regexp, String)` 方法并不会修改原来的字符串，它会先拷贝原字符串，对其修改之后返回。如果 *String* 对象包含的内容很多，拷贝可能比较消耗资源，因此 Ruby 也提供了直接对原字符串生效的方法 `String.sub!(Regexp, String)`。

```
string = "2010-12-20 2011-02-14"
string.sub(/\d{4}-\d{2}-\d{2}/, 'date')
puts string
date 2011-02-14
```

如果表达式在 *String* 中能多次匹配，`String.sub()` 只会对第一次匹配进行替换；如果需要对所有匹配都进行替换，则需要使用 `String.gsub(Regexp, String)`。下面来看这个方法。

### 15.3.9 String.gsub(Regexp, String)

这个方法在使用上完全等价于 `String.sub(Regexp, String)`，唯一不同的是它会对每次匹配都进行替换。

```
puts "2010-12-20 2011-02-14".gsub(/\d{4}-\d{2}-\d{2}/, '\2/\3/\1')
12/20/2010 02/14/2011
```

相应的，`String.gsub!(Regexp, String)` 也会修改原来的字符串。

```
string = "2010-12-20 2011-02-14"
string.gsub!(/\d{4}-\d{2}-\d{2}/, '\2/\3/\1')
puts string
12/20/2010 02/14/2011
```

## 15.4 常用操作示例

### 15.4.1 验证

简单验证

```
#不要忘记在正则表达式首尾加上\A 和\z
("2010-12-20" =~ /\A\d{4}-\d{2}-\d{2}\z/) != nil    #=> True

#也可以使用 String.index(Regexp)
#同样要在正则表达式首尾加上\A 和\z
"2010-12-20".index(/\A\d{4}-\d{2}-\d{2}\z/) != nil    #=> True
```

## 15.4.2 提取

### 逐步提取

```
regex = /(\d{4})-(\d{2})-(\d{2})/
text = "2010-12-20 2011-02-14"
data = regex.match(text)
while data != nil
  print "date:", data[0], " year:", data[1], " month:", data[2], " day:", data[3], "\n"
  data = regex.match(text, data.end(0))
end
```

```
date: 2010-12-20, year: 2010, month: 12, day: 20
```

```
date: 2011-02-14, year: 2011, month: 02, day: 14
```

### 一次性提取

```
regex = /((\d{4})-(\d{2})-(\d{2}))/
text = "2010-12-20 2011-02-14"
all = text.scan(regex)
for one in all do
  print "date:", one[0], " year:", one[1], " month:", one[2], " day:", one[3], "\n"
end
```

```
date: 2010-12-20, year: 2010, month: 12, day: 20
```

```
date: 2011-02-14, year: 2011, month: 02, day: 14
```

## 15.4.3 替换

### 简单替换

```
puts "2010-12-20 2011-02-14".gsub(/(\d{4})-(\d{2})-(\d{2})/, '\2/\3/\1')
12/20/2010 02/14/2011
```

### 在 replacement 中使用整个表达式匹配的文本

```
#必须给整个表达式添加捕获型括号
regex = /(\d+\.\d{0,2})/
replacement = '$\0'
"the price is 12.99".sub(regex, replacement)
the price is $12.99
```

## 15.4.4 切分

```
for word in "one\ttwo\n three".split(/\s+/) do
  puts word
end

one two three
```

## 15.5 Ruby 1.9 的新变化

Ruby 1.8 的正则表达式还不够完善,许多常用的特性不能支持,而 Ruby 1.9 采用了全新的 Oniguruma (鬼车) 引擎,正则表达式支持则强了很多。下面列出 Ruby 1.9 中正则表达式的主要变化,更详细的信息请参考网址 <http://oniguruma.rubyforge.org/oniguruma/>。

### 新增了合并多个表达式的方法

`Regexp.union(Regexp, Regexp, ...)`

这个方法返回的仍然是一个 `Regexp`, 它“合拢”了参数中所有的 `Regexp`, 任何一段文本, 只要能被其中某个表达式匹配, 就能被 `Regexp` 匹配。

```
"a" =~ Regexp.union(/^d+$/, /^[a-z]+$/) # => 0
"0" =~ Regexp.union(/^d+$/, /^[a-z]+$/) # => 0
```

关于这个方法, 有一个常见的误区: 在数据验证时, 想当然地将各条件对应的表达式用 `Regexp.union()` “归并”, 但如果这些条件必须同时满足, 就很可能出错。

比如需要验证的字符串有这样两条要求: 第一, 全部由数字或字母组成; 第二, 长度在 6~12 之间。如果想当然地用 `Regexp.union()`, 就会出错。

```
"0123" =~ Regexp.union(/^[a-zA-Z0-9]*$/, /^.{6,12}$/) # => 0
```

这是因为 `Regexp.union()` “归并”的表达式并不要求“同时生效”, 只要其中有一条能匹配, 则匹配成功。此处字符串 `0123` 可以由表达式 `^[a-zA-Z0-9]*$` 匹配, 所以返回结果 0。

### 新增了逆序环视

```
"ab" =~ /ab(?<=ab)/ # => 0
"cd" =~ /cd(?<!ab|cde)/ # => 0
```

### 新增了忽略优先量词

```
puts /\.*/?.match('"string 1" and "string 2"')
"string 1"

puts /\.*/?.match('"string 1" and "string 2"')
"string 1" and "string 2"
```

### 新增了字符组中的集合运算操作符&&

```
"a" =~ /[a-z&&[aeiou]]/ # => nil
"b" =~ /[a-z&&[aeiou]]/ # => 0
```

### 新增了指定 Unicode 码值的表示法

如果指定了多字节编码，十进制和十六进制的数字序列可以用来表示单个多字节字符，因此可以在字符组中指定多字节字符的范围。

```
"你" =~ /[\\u{4e00}-\\u{9fff}]/u # => 0
"a"  =~ /[\\u{4e00}-\\u{9fff}]/u # => nil
```

如果不知道自己使用的 Ruby 的版本，或者不确认其使用了 Oniguruma 引擎，则可以使用下面的函数来判断。

```
def oniguruma?
  return true if RUBY_VERSION >= "1.9.0"

  if defined?(Regexp::ENGINE) # Is ENGINE defined
    if Regexp::ENGINE.include?('Oniguruma')
      return true           # Some version of Oniguruma
    else
      return false          # Pre-Oniguruma engine
    end
  end
end
```

如果希望无缝兼容 Ruby 1.8 和 1.9，不妨自己包装一个方法，根据要处理的字符串来生成对应编码的 Regexp。

```
def get_regex(pattern, str, options=0)
  if str.methods.include?(:encoding) # Ruby 1.8 compatibility
    return Regexp.new(pattern.encode(encoding),options)
  else
    return Regexp.new(pattern, options)
  end
end
```



## 第 16 章 Linux/UNIX

本章介绍 Linux 下常用工具（vi、grep、awk、sed）中正则表达式的应用。虽然大多数工具在 Linux/UNIX/FreeBSD 中都存在，但版本和用法都存在若干差异。本章的讲解以 Linux 为准，如果你使用其他平台，可以参考本章的内容，具体细节则需要参考对应的文档。

### 16.1 POSIX

如果你用过 vi、grep、awk、sed 之类的工具，或许会记得：这些工具虽然支持正则表达式，但语法却很不一样，按照通常习惯的办法写的正则表达式往往不是无法识别（`(abc)` 竟然可能报错）就是匹配错误（`\d` 无法匹配数字字符）。而且，这些工具自身之间也存在差异，同样的结构，有时需要转义有时不需要转义。这究竟是因为什么呢？

原因在于，Linux/UNIX 下的工具大多采用 POSIX 规范，而 POSIX 规范又可分为两种流派（flavor）。之前只是简单介绍了 POSIX 字符组，现在详细讲解 POSIX 规范。

#### 16.1.1 POSIX 规范

我们之前介绍的正则表达式的记法，其实都源于 Perl，实际上，正则表达式从 Perl 衍生出一个显赫的流派，叫做 PCRE（Perl Compatible Regular Expression），`\d`、`\w`、`\s` 之类的字符组简记法，是它的特征。但是在 PCRE 之外，正则表达式还有其他流派，比如 POSIX 规范的正则表达式。

POSIX 的全称是 Portable Operating System Interface for uniX，它由一系列规范构成，定义了 UNIX 操作系统应当支持的功能，所以“POSIX 规范的正则表达式”其实只是“关于正则表达式的 POSIX 规范”，它定义了 BRE（Basic Regular Expression，基本型正则表达式）和 ERE（Extended Regular Expression，扩展型正则表达式）两大流派。在兼容 POSIX 的 UNIX 系统上，grep 和 egrep 之类的工具都遵循 POSIX 规范，一些数据库系统中的正则表达式也符合 POSIX 规范。

##### 16.1.1.1 BRE

grep、vi、sed 属于 BRE 这一派，它的语法看起来比较奇怪，元字符 `[`、`]`、`{`、`}` 必须转义之后才具有特殊含义，比如 `(a)b` 只能匹配 `(a)b` 而不是 `ab`；`a{1,2}` 只能匹配字符串 `a{1,2}`，`a\{1,2\}` 才能匹配字符串 `a` 或者 `aa`。

之所以这么麻烦，是因为正则表达式的许多功能是逐步出现的，一开始这些元字符并没有特

殊的含义。为保证向后兼容，就只能使用转义。基于同样的原因，BRE 不支持+和?量词，也不支持多选结构(...)和反向引用\1、\2…。

缺少了这些功能，BRE 使用起来确实不够方便，好在今天纯粹的 BRE 已经很少见了，毕竟大家已经认为正则表达式“理所应当”支持多选结构和反向引用等功能。所以虽然 vi 属于 BRE 流派，仍然提供了这些功能。而且，GNU 也对 BRE 做了扩展，支持+、?、|，只是使用时必须写成\+、\?、\|，而且也支持\1、\2之类的反向引用。这样，GNU 的 grep 等工具虽然名义上属于 BRE 流派，但更确切的名称是 GNU BRE。

16.1.1.2 ERE

egrep、awk 则属于 ERE 这一派，虽然从名字上来看，BRE 是“基本”而 ERE 是“扩展”，但 ERE 不要求兼容 BRE 的语法，而是自成一体。因此其中的元字符不用转义（在元字符之前添加反斜线会取消其特殊含义），所以(ab|cd)就可以匹配字符串 ab 或者 cd，量词+、?、{n,m}可以直接使用。ERE 并没有明确规定支持反向引用，但是不少工具都支持\1、\2之类的反向引用。

GNU 出品的 egrep 等工具就属于 ERE 流派（更准确的名字是 GNU ERE），但因为 GNU 已经对 BRE 做了不少扩展，所谓的 GNU ERE 其实只是个说法，它的功能 GNU BRE 都有了，GNU ERE 的主要区别在于元字符不需要转义。

表 16-1 简要说明了几种 POSIX 流派的区别<sup>1</sup>（其实，现在的 BRE 和 ERE 在功能上并没有什么区别，主要的差异是在元字符的转义上）。

表 16-1 几种 POSIX 流派的区别

流派	说明	工具
BRE	[、]、{、}都必须转义使用，不支持+、?、	grep、sed、vi（但 vi 支持多选结构和反向引用）
GNU BRE	[、]、{、}、+、?、 都必须转义使用	GNU grep、GNU sed
ERE	元字符不必转义，+、?、[、]、{、}、 可以直接使用，对\1、\2之类的支持不确定	egrep、awk
GNU ERE	元字符不必转义，+、?、[、]、{、}、 可以直接使用，支持\1、\2之类	grep-E、GNU awk

16.1.2 POSIX 字符组

POSIX 流派与 PCRE 流派的一大区别是字符组。

通常（也就是 PCRE 流派的做法）我们把[...]称为“字符组”；但在 POSIX 中，它们叫做“方括号表达式”（bracket expression），其作用与规则同常说的“字符组”完全一样：在方括号中列

<sup>1</sup> 关于 ERE 和 BRE 的详细规范，可以参考 [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html)。

出可以匹配的所有字符，比如[abc]，也支持横线表示范围，比如[a-z]，如果在方括号里的开头写脱字符^，则表示排除型字符组，比如[^0-9]。但是，POSIX 中的“方括号表达式”都不支持使用\uxxxx 的形式表示 Unicode 字符。也就是说，遵循 POSIX 规范的正则表达式无法用[\u4e00-\u9fff]来匹配“任意一个中文字符”。

POSIX 自己也定义了字符组，准确的名称是“POSIX 字符组”(POSIX character class)，这是一些事先定义好的组合，基本等价于我们常说的“字符组简记法”。比如[:digit:]就是一个 POSIX 字符组，它等价于[0-9]，[:alpha:]也是 POSIX 字符组，它等价于[a-zA-Z]。另一方面，PCRE 中常见的\d、\w、\s 在 POSIX 中并没有定义。

POSIX 字符组不能直接使用，只能用于方括号表达式内部，类似[[:digit:]]和[[:alpha:]]；POSIX 字符组也可以与其他方括号表达式合用，比如[a-zA-F[:digit:]]；[^[:alpha:]]匹配一个非字母字符，[^a-zA-F[:digit:]]匹配一个十六进制数字之外的字符。

表 16-2 详细介绍了常用的 POSIX 字符组（请注意，POSIX 字符组只匹配 ASCII 字符）。

表 16-2 常用的 POSIX 字符组

POSIX 字符组	Perl 记法	ASCII 编码能匹配字符	说明
[[:alnum:]]	类似 \w，但不能匹配下划线	[A-Za-z0-9]	数字和字母字符
[[:alpha:]]		[A-Za-z]	字母字符
[[:ASCII:]]		[\x00-\x7F]	ASCII 字符
[[:blank:]]		[ \t]	空格和制表符
[[:cntrl:]]		[\x00-\x1F\x7F]	控制字符
[[:digit:]]	\d	[0-9]	数字字符
[[:graph:]]		[\x21-\x7E]	可见字符
[[:lower:]]		[a-z]	小写字母字符
[[:print:]]		[\x20-\x7E]	可见字符和空白字符
[[:punct:]]		[!\"#\$%&'()*+,-./:;<=>?@^_`{ }~]	标点符号
[[:space:]]	\s	[ \t\r\n\v\f]	空白字符
[[:upper:]]		[A-Z]	大写字母字符
[[:xdigit:]]		[A-Fa-f0-9]	十六进制数字字符

了解了 POSIX、BRE、ERE 这些概念之后，下面简要介绍 Linux/UNIX 下常用工具中的正则表达式。

## 16.2 vi

vi/vim（vim 是 vi iMproved 的缩写，以下统称 vi）是 Linux/UNIX 下常用的文本编辑器，它提供了各种正则表达式功能。vi 的正则表达式属于 BRE 流派，但是它也支持几乎所有的正则表达式功能，只是不少元字符都需要转义。如果使用 Very Magic 模式（后面会讲到），也可以不用那么麻烦的转义。

首先介绍查找所用的各种记法。在 vi 里进行查找，首先要切换到控制模式（按 Esc 键），然后输入 /*pattern*/（结尾的 / 也可以忽略），其中 *pattern* 就是用于查找的模式字符串。在控制模式下输入 /cat/ 再回车，就查找当前编辑的文本中所有的字符串 cat（注意不是单词 cat）。

### 16.2.1 字符组及简记法

vi 的字符组遵循 POSIX 规范，所以 POSIX 中的各种字符组和记法都可以直接在 vi 中使用。如果在开方括号之前写上 \，则这个字符组不但能匹配其中列出的任何字符，还可以匹配行终止符，比如 \[abc] 就不但能匹配字符 a、b、c，还能匹配换行符。

除去 POSIX 规范的字符组，vi 也提供了类似 Perl 的字符组简记法，虽然 vi 在正确设定编码后<sup>1</sup>能够识别中文等多字节字符（不会把一个中文字符当成“几个字符”来处理），但其对“数字字符”、“空白字符”和“单词字符”的识别仍然局限于 ASCII 字符，详见表 16-3。

表 16-3 vi 中的字符组简记法

字符组简记法	说明	字符组简记法	说明
.	除行终止符之外的任何字符	\.	匹配任何字符（包括行终止符）
\s	ASCII 编码中的空白字符（但不包括行终止符）	\S	\s 无法匹配的字符
\d	[0-9]	\D	\d 无法匹配的字符
\x	[0-9a-fA-F]	\X	\x 无法匹配的字符
\o	[0-7]	\O	\o 无法匹配的字符
\h	[a-zA-Z_]	\H	\h 无法匹配的字符
\p	可以打印的字符（非控制字符）	\P	除数字字符之外，所有的可打印字符
\w	[0-9a-zA-Z_]	\W	\w 无法匹配的字符
\a	[a-zA-Z]	\A	\a 无法匹配的字符

<sup>1</sup> 在 vi 中，切换到控制模式下，输入 :set enc=xxx 可以设定编码，其中 xxx 可以是 utf8 或者 gbk，注意在 set 之前有一个冒号。

( 续表 )

字符组简记法	说明	字符组简记法	说明
<code>\l</code>	<code>[a-z]</code>	<code>\L</code>	非小写字母字符
<code>\u</code>	<code>[A-Z]</code>	<code>\U</code>	非大写字母字符

补充一点，表中列出的字符组简记法是不能在方括号内部使用的。比如匹配一个十六进制数字字符，可以用`\x`，也可以用`[[[:digit:]]a-fA-F]`，但`[\da-fA-F]`无法匹配`[0-9]`。

要值得注意的是，vi 中没有单行模式，所以点号不能匹配行终止符，如果需要匹配包括行终止符在内的任何字符，可以使用`\_.`。

16.2.2 量词

vi 中的量词比较奇怪，大多需要转义，部分通行量词的记法也有区别（主要是`{`出现之后，`}`不需要转义），刚开始接触可能不习惯。表 16-4 列出了 vi 中的量词。

表 16-4 vi 中的量词

量词	等价的 PCRE 量词
<code>*</code>	<code>*</code>
<code>\+</code>	<code>+</code>
<code>\=</code>	<code>?</code>
<code>\{n,m}</code>	<code>{n,m}</code>
<code>\{n}</code>	<code>{n}</code>
<code>\{n,}</code>	<code>{n,}</code>
<code>\{,m}</code>	<code>{0,m}</code>
<code>\{-}</code>	<code>*?</code>
<code>\{-n,m}</code>	<code>{n,m}?</code>
<code>\{-n,}</code>	<code>{n,}?</code>
<code>\{-,m}</code>	<code>{,m}?</code>

vi 没有现成的忽略优先量词`+?`和`??`，但是可以用`\{-1, }`和`\{-, 1}`表示。

16.2.3 多选结构和捕获分组

前面讲过 vi 的正则表达式属于 BRE 流派，所以`[`、`|`、`)`都需要转义。通常我们所见的`(...|...)`，要写作`\(...\|...\)`。

捕获分组中的括号也需要转义，所以`(...)`要写作`\(...\)`。

## 16.2.4 环视

vi 也提供了环视功能，但是它的环视写法有所不同。vi 的环视与 PCRE 的环视对照如表 16-5 所示。

表 16-5 vi 的环视与 PCRE 的环视对照举例

vi	PCRE
<code>\@=</code>	<code>?=</code>
<code>\@!</code>	<code>?!</code>
<code>\@&lt;=</code>	<code>?&lt;=</code>
<code>\@&lt;!</code>	<code>?&lt;!</code>

同时环视的写法也与常见的不同，PCRE 将环视使用的表达式写在环视的后方（右方），而 vi 中必须将环视使用的表达式写在环视的左方，表 16-6 举了两个例子。

表 16-6 vi 的环视与 PCRE 的环视对照举例

vi	PCRE
<code>ab\((cd\)\)\@=</code>	<code>ab(?:=cd)</code>
<code>\(foo\)\@&lt;!bar</code>	<code>(?:&lt;!foo)bar</code>

vi 的环视功能没有任何限制，可以随意使用量词，所以可以很方便地用它查找“包含/不包含某个表达式”的行。

查找包含 tom 或者 jerry 的行。

```
^\(.*(tom|jerry)\)\@=
```

查找不包含 tom 也不包含 jerry 的行。

```
^\(.*(tom|jerry)\)\@!
```

## 16.2.5 锚点和单词边界

vi 提供了常见的匹配行起始位置和结束位置的 `^` 和 `$`，但没有提供“通用”的单词边界 `\b`，而是提供了分别匹配单词起始和结束位置的 `\<` 和 `\>`。详见表 16-7。

表 16-7 vi 中的锚点和单词边界

分界符	说明
<code>^</code>	匹配行的起始位置
<code>\$</code>	匹配行的结束位置
<code>\&lt;</code>	单词开始位置，左边不能是 <code>[0-9a-zA-Z_]</code> ，右边必须是 <code>[0-9a-zA-Z_]</code>
<code>\&gt;</code>	单词结束位置，右边不能是 <code>[0-9a-zA-Z_]</code> ，左边必须是 <code>[0-9a-zA-Z_]</code>

下面举几个例子。

匹配行开始的 cat

```
^cat
```

匹配行结束的 cat

```
cat$
```

匹配单词开始的 cat

```
\<cat
```

匹配单词结束的 cat

```
cat\>
```

回到开头提到的查找 cat 的例子，如果我们要查找的是“单词”cat，就应该写作 `\<cat\>`。

## 16.2.6 替换操作的特殊字符

vi 中的替换与查找类似，都是先切换到控制模式（按 Esc 键），然后输入 `:x,ys/pattern/replacement/`。其中 *x* 和 *y* 用来限制替换操作所在的行数（从第 *x* 行开始，到第 *y* 行结束，如果不指定 *x* 和 *y*，则只在当前行进行替换；如果只写一个%，则对所有行进行替换），~~s~~ 表示“替换”操作，*pattern* 是用来查找的模式字符串，*replacement* 则是要替换成的字符串，当然也不要忘记最前面的冒号。

在当前行将单词 tom 替换为 jerry。

```
:s/\<tom\>/jerry/
```

在第 1 行到第 10 行将单词 tom 替换为 jerry。

```
:1,10s/\<tom\>/jerry/
```

在所有行将单词 tom 替换为 jerry。

```
:%s/\<tom\>/jerry/
```

为简明起见，下面的替换命令中不再写出 *x,y*，只保留 `s/pattern/replacement/`，请读者注意。

在默认情况下，替换只会对第一次出现的进行，如果在某行中多次出现了正则表达式能够匹配的文本，需要将它们全部替换，可以使用 `s/pattern/replacement/g`，最后的 *g* 表示“全部替换”模式。

在所有行中将出现的所有单词 tom 替换为 jerry。

```
:%s/\<tom\>/jerry/g
```

在结尾的/之后，除了 g，还有其他模式，如表 16-8 所示。

表 16-8 vi 中能使用的模式

模式	说明
g	全局性替换（尽可能多地替换）
i	查找时不区分大小写
c	在替换前需要提示确认

下面再给出一些常见操作的写法。

在所有行中将所有出现的单词 tom（不区分大小写）替换为 jerry，并在替换前确认：

```
:%s/\<tom\>/jerry/igc
```

删除所有行开头的空白

```
s/^s\+//
```

删除所有行结尾的空白

```
s/s\+$//
```

删除所有只包含空白字符的行中的空白字符

```
%s/^s\+$//g
```

行开头添加#（Python 和 Ruby 注释）

```
%s/^/#/
```

删除行末的注释#

```
s/#s\+$//
```

将连续空行合并为一个空行

```
s/^\_[:space:]\+$/g
```

删除 HTML tag

```
s/<[^>\+>//g
```

删除所有包含 tom 或 jerry 的行

```
s/^\(.*(tom|jerry)\)\@=.*$/i
```

删除所有//注释

```
s/\/\/.*/
```



16.2.7 replacement 中的特殊变量

表 16-9 中列出了 replacement 中可以使用的特殊变量。

表 16-9 replacement 中可以使用的特殊变量

变量	说明
&	表达式匹配的全部文本
\0	表达式匹配的全部文本
\1	编号为 1 的捕获分组匹配的文本
\2	编号为 2 的捕获分组匹配的文本
\L	将之后的所有字符转换为小写
\U	将之后的所有字符转换为大写
\E	\U 和 \L 的作用结束
\e	\U 和 \L 的作用结束
\l	将此后一个字符转换为小写
\u	将此后一个字符转换为大写

用<b>标签标注 word。

```
s/\<word\>/<b>&</b>/g
```

删除重复单词。

```
s/\<(\w+)\>\s+\<\1\>/\1/g
```

找出所有形式的 tom 或者 jerry，全部转换为小写。

```
s/\< \(tom\|jerry\) \>/\L&/gi
```

将所有的 tom 或 jerry 统一为首字母大写的形式。

```
s/\< \(tom\|jerry\) \>/\u&/g
```

16.2.8 补充

习惯 PCRE 正则表达式的用户可能不习惯 vi 中频繁出现的转义。如果是这样，可以使用 vi 提供的 Very Magic 模式，在这种模式下，常用的 `$`、`.`、`*`、`(`、`)`、`|` 都不用转义，只有 `{` 和 `}` 需要转义。

Very Magic 模式的使用非常简单，在查找字符串的最开始写上 `\v` 就可以。表 16-10 给出了两个例子。

表 16-10 Very Magic 模式使用举例

普通模式	Very Magic 模式
\(...\ ...\)	\v(... ...)
\(...\)	\v(...)

在 vi 中输入 :h pattern 可以查询到 pattern 的详细信息，也可以参考这两份文档。

[http://en.wikipedia.org/wiki/Vim\\_\(text\\_editor\)](http://en.wikipedia.org/wiki/Vim_(text_editor))

<http://vimdoc.sourceforge.net/>

## 16.3 grep

grep 是 Linux/UNIX 下常用的工具，作用是从文件或规范输入中筛选出包含某个正则表达式能匹配文本的行，将其传送至标注输出。grep 名字来自它所模拟的 ed 编辑器的一系列操作的名称缩写：Global Regular Expression Print，也有种说法是 Generalized Regular Expression Parser，最早出现在 1973 年的 UNIX Version 4 中。目前大多数 Linux/UNIX 中都提供了 grep（Windows 下也有 grep，但是需要自己下载），其中以 GNU grep 最为流行，下面的讲解就以 GNU grep 为准。

### 16.3.1 基本用法

grep 的基本用法是：

grep [options] PATTERN [FILES]

其中的 options 在本节最后介绍，PATTERN 是需要指定的正则表达式，FILES 是需要检索的文件，如果不指定 FILES，则默认从标准输入（stdin）读取文本，下面举一个最简单的例子。

从 sample.txt 中找出所有包含 word 的行。

grep word sample.txt

为了避免正则表达式被 Shell 错误解释（正则表达式中的 \*、[、] 之类可能先被 Shell 展开），保险的做法（也是好的习惯）是给正则表达式加上单引号'。

grep 'word' sample.txt

如果在其中出现了变量插值，则应当使用双引号"。

grep "\$JAVA\_HOME" sample.txt

### 16.3.2 字符组

grep 的字符组遵循 POSIX 规范，所以 POSIX 的字符组和简记法都可以直接在 grep 中使用。grep 中的点号可以匹配除换行符之外的任何字符，因为 grep 是按行处理的，每次处理一行，所

以不需要关心点号能否匹配换行符。

从 sample.txt 中找出所有这样的行：包含字符 a 或 b

```
grep '[ab]' sample.txt
```

从 sample.txt 中找出所有这样的行：包含除 a 和 b 之外的任何字符

```
grep '[^ab]' sample.txt
```

从 sample.txt 中找出所有这样的行：包含 cat、cut、cit 之类文本的行

```
grep 'a.b' sample.txt
```

### 16.3.3 锚点和单词边界

grep 中能使用的锚点有 4 个：`^`、`$`、`\<`、`\>`。`^`用来匹配行的开始位置，`$`用来匹配行的结束位置（因为 grep 是面向行的处理工具，每次处理一行，所以不存在`^`和`$`匹配文本内部的行起始/结束位置的问题，也就不需要多行模式）。grep 不支持常见的`\b`，而是用`\<`和`\>`来匹配单词的起始位置和结束位置，它们所识别的“单词字符”是`[:alnum:]`能匹配的字符，也就是`[0-9a-zA-Z]`，注意其中不包含下划线。

从 sample.txt 中找出所有以 cat 开头的行

```
grep '^cat' sample.txt
```

从 sample.txt 中找出所有以 cat 结尾的行

```
grep 'cat$' sample.txt
```

从 sample.txt 中找出所有包含单词 cat 的行

```
grep '\<cat\>' sample.txt
```

### 16.3.4 量词

grep 中能使用的量词有`*`、`+`、`?`、`{n, m}`，其中只有`*`可以直接使用，其他量词都必须转义：`+`必须写成`\+`，`?`必须写成`\?`，`{n, m}`要写成`\{n, m\}`。

从 sample.txt 中找出所有包含 a、ab、abb、abbb...的行

```
grep 'ab*' sample.txt
```

从 sample.txt 中找出所有包含 ab、abb、abbb...的行

```
grep 'ab\+' sample.txt
```

从 sample.txt 中找出所有包含 a 或者 ab 的行

```
grep 'ab?' sample.txt
```

从 sample.txt 中找出所有包含 ab、abb、abbb 的行

```
grep 'ab\{1,3\}' sample.txt
```

### 16.3.5 多选结构和捕获分组

grep 支持多选结构，但是必须转义 `(`、`|`、`)` 三个字符，多选结构必须写成 `\(...\|...\)`，捕获分组的记法是 `\1`、`\2`……。

查找 sample.txt 中所有包含 tom 或者 jerry 的行。

```
grep '\(tom\|jerry\)' sample.txt
```

查找 sample.txt 中所有包含重复单词的行。

```
grep '\(<[[:alpha:]]\+\>\)[[:space:]]\+\<\1\>' sample.txt
```

### 16.3.6 options

grep 中可以使用 *options* 来指定匹配规则（类似匹配模式），常用的 *option* 介绍如下。

**-i 不区分大小写**

查找各种形式的 cat，包括 CAT、cat、Cat 等

```
grep -i 'cat' sample.txt
```

**-v 不匹配**

这是 grep 中特别有用的功能，正则表达式要“排除若干字符串”通常是非常麻烦的，但是有了 `grep -v` 就非常方便了，它会选出正则表达式不能匹配的所有行。

查找所有不包含 cat 的行

```
grep -i 'cat' sample.txt
```

查找所有不包含 404 和 302 的行（在过滤日志时特别有用）

```
grep -v '\(404\|302\)' sample.txt
```

查找所有不包含 jpg、png、gif 的行（在过滤日志时特别有用）

```
grep -v '\(jpg\|png\|gif\)' sample.txt
```

**-w 只匹配单词**

查找只包含单词 cat 的行，等价于 `grep '\<cat\>' sample.txt`

```
grep -w 'cat' sample.txt
```

**-n 显示行数**

找出所有包含 `cat` 的行，同时标明行数

```
grep -n 'cat' sample.txt
```

**-r 递归查找**

这也是一个特别有用的功能，它可以对目录下所有文件进行查找。

找出当前目录下所有文件中包含 `cat` 的行

```
grep -r 'cat' .
```

**-A *n* 显示匹配行及之后的 *n* 行**

找出所有包含 `cat` 的行，并显示每行之后的 3 行

```
grep -A 3 'cat'
```

**-B *n* 显示匹配行及之前的 *n* 行**

找出所有包含 `cat` 的行，并显示每行之前的 3 行

```
grep -B 3 'cat'
```

**-C *n* 显示匹配行及之前之后各 *n* 行（等价于 `-A n -B n`）**

找出所有包含 `cat` 的行，并显示每行之前和之后的 3 行

```
grep -C 3 'cat'
```

**-P 使用 Perl 风格的表达式**

如果你很熟悉 Perl 风格的正则表达式，而不习惯 `grep` 的规定，开启这个选项可以直接使用 Perl 风格的正则表达式，比如像下面这样查找单词 `cat`。

```
grep -P '\bcat\b'
```

### 16.3.7 egrep 和 fgrep

`egrep` 和 `grep` 的名字不同，但并没有本质区别。虽然 `egrep` 属于 ERE 流派，而 `grep` 属于 BRE 流派，但 GNU `egrep` 和 GNU `grep` 其实只有写法的差别，功能上并无不同，实际上，在命令行下输入 `grep -E`，就等价于 `egrep`。使用 `egrep` 的理由或许在于，因为它属于 GNU ERE 流派，元字符几乎不需要转义就可直接使用，写、读起来都要更简单一些。

`fgrep` 也是 `grep` 家族中的成员，`f` 代表 fast，`fgrep` 也确实比 `grep` 要快很多，因为它完全取消了正则表达式的功能，只进行最简单的字符串查找。在 `fgrep` 中查找 `end$`，并不是找出行末尾为 `end` 的行，而是包含字符串 `end$` 的行。

### 16.3.8 补充

`grep` 负责选出能匹配成功的整行文本，真正匹配到的内容反而不好识别了。使用 `grep` 时最

好加上一个参数 `grep --color`，这样 `grep` 会把能匹配到的文本用不同颜色标注出来，方便识别。如果你愿意，不妨做一个 `alias`，省去每次明确指定：

```
alias grep='grep --color'
```

关于 `grep` 的详细文档，可以输入 `info grep` 获得，也可以参考下面的文档：

<http://www.gnu.org/software/grep/manual/>

<http://en.wikipedia.org/wiki/Grep>

## 16.4 awk

`awk` 也是 Linux/UNIX 下常用的工具，它在 20 世纪 70 年代诞生于贝尔实验室，得名自三位作者 Alfred Aho、Peter Weinberger、Brian Kernighan。`awk` 的主要作用是对数据重新组织和统计，它功能非常多，语法也很复杂，如果使用得当完全可以实现非常复杂的逻辑。在本节，我们只介绍与正则表达式相关的内容。`awk` 有许多版本，本章以 `gawk`（GNU `awk`）为准。

### 16.4.1 基本用法

`awk` 的基本使用方法是：

```
awk 'program' input-file1 input-file2 ...
```

其中 *program* 是对数据处理程序代码，其基本形式为：

```
condition { action }
```

*condition* 是判断条件，*action* 是条件满足时所执行的操作，比如下面这条命令的意思就是，如果 `sample.txt` 中的某一行可以由正则表达式 `J` 匹配（也就是说，包含大写字母 `J`），就输出这一整行（`$0` 的含义在下面讲解）。

```
awk '/J/ {print $0}' sample.txt
```

为方便讲解，假设 `sample.txt` 的内容如下：

```
Jan 13 25 15 115
Feb 15 32 24 226
Mar 15 24 34 228
Apr 31 52 63 420
May 16 34 29 208
Jun 31 42 75 492
Jul 24 34 67 436
```

上面命令运行的结果就是：

```
Jan 13 25 15 115
Jun 31 42 75 492
Jul 24 34 67 436
```

如果仅仅这样使用，`awk` 和 `grep` 并没有太大的区别，但 `awk` 的功能并不只这些。`awk` 提供了一系列变量 `$1`、`$2`、`$3`... 对应每一行的各个字段（在默认状态下，字段之间以空白分隔，也可以通过 `FS` 指定分隔），`$0` 代表整行，`NF` 代表该行的字段数目，因此 `$NF` 表示最后一个字段。

如果把上面的语句改一改，就可以同时实现对字段的过滤。

```
awk '/J/ {print $1, $2, $NF}' sample.txt
```

```
Jan 13 115
Jun 31 492
Jul 24 436
```

我们还可以在 `{action}` 中写上一些处理逻辑。

```
awk '{if( $NF > 400) print $0}' sample.txt
```

```
Apr 31 52 63 420
Jun 31 42 75 492
Jul 24 34 67 436
```

下面介绍 `awk` 中的正则表达式。`awk` 中的正则表达式虽然也遵循 POSIX 规范，但属于 ERE 流派，所以很多元字符不需要麻烦的转义。

## 16.4.2 字符组及简记法

`awk` 的字符组遵循 POSIX 规范，所以 POSIX 的字符组和简记法都可以直接在 `awk` 中使用。而且，`awk` 中的点号可以匹配包括换行符在内的任何字符，不过一般来说 `awk` 是按行处理的，所以这一点并不重要。

从 `sample.txt` 中找出所有包含 `a` 或者 `b` 的行。

```
awk '/[ab]/ {print $0}' sample.txt
```

从 `sample.txt` 中找出所有包含数字字符的行。

```
awk '/[0-9]/ {print $0}' sample.txt
awk '/[[:digit:]]/ {print $0}' sample.txt
```

例外的是，GNU `awk` 支持两个 PCRE 的字符组 `\w` 和 `\W`，其中 `\w` 等价于 `[[:alnum:]]`，`\W` 则匹配 `\w` 所不能匹配的所有字符。

### 16.4.3 锚点和单词边界

awk 中能使用的锚点有下面几个：`^`、`$`、`\<`、`\>`、`\y`、`\b`。`^`用来匹配行的开始位置，`$`用来匹配行的结束位置（因为 awk 是面向行的处理工具，每次处理一行，所以不存在`^`和`$`匹配文本内部的行起始和结束位置的问题，也就不需要多行模式）。

awk 识别的“单词字符”是`[[:alnum:]]`能匹配的字符，也就是`[0-9a-zA-Z]`，用`\<`和`\>`匹配单词的起始和结束位置，`\y`匹配单词的起始或者结束位置，`\b`匹配`\y`无法匹配的位置。

从 sample.txt 中找出所有以 cat 开头的行。

```
awk '/^cat/ {print $0}' sample.txt
```

从 sample.txt 中找出所有以 cat 结尾的行。

```
awk '/cat$/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含以 cat 开头单词的行。

```
awk '/\<cat/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含以 cat 结尾单词的行。

```
awk '/cat\>/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含单词 cat 的行。

```
awk '/ \<cat\>/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含 cat，但不是单词 cat 的行（比如 truncate）。

```
awk '/\Bcat\b/ {print $0}' sample.txt
```

### 16.4.4 量词

awk 中能使用的量词有`*`、`+`、`?`、`{n,m}`、`{n,}`、`{n}`，因为 awk 属于 ERE 流派，所以这些量词可以直接使用，不用转义。

从 sample.txt 中找出所有包含 a、ab、abb、abbb...的行。

```
awk '/ ab*/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含 ab、abb、abbb...的行。

```
awk '/ ab+/ {print $0}' sample.txt
```

从 sample.txt 中找出所有包含 a 或者 ab 的行。

```
awk '/ ab?/ {print $0}' sample.txt
```



从 sample.txt 中找出所有包含 ab、abb、abbb 的行。

```
awk '/ ab{1,3}/ {print $0}' sample.txt
```

### 16.4.5 多选结构

awk 中的多选结构也是 `(...|...)` 的形式，可以直接使用，不需要转义。

从 sample.txt 中找出包含 tom 或者 jerry 的行。

```
awk '/(tom|jerry)/ {print $0}' sample.txt
```

### 16.4.6 补充

awk 不支持捕获分组和反向引用，不支持忽略优先量词，也无法使用匹配模式进行不区分大小写查找。

关于 awk 的详细文档，可以输入 `info awk` 获得，也可以参考下面的文档：

<http://www.gnu.org/manual/gawk/gawk.html>

## 16.5 sed

sed 也诞生于贝尔实验室，是 Linux/UNIX 下常用的流编辑工具。sed 得名自 stream editor，它主要用来进行文本的替换等操作，因为 sed 进行流式处理（在默认模式下 sed 每次读入一行，去掉换行符，处理，再加上换行符，输出），所以速度非常快。sed 有许多版本，这里讲解的以 GNU sed 为准（选择它的理由是因为它比较标准，也很方便，GNU sed 提供了不区分大小写的 i 模式，其他的一些 sed 则没有提供）。前面讲过，sed 属于 BRE 流派，所以许多元字符之前需要加反斜线来转义。

### 16.5.1 基本用法

sed 最主要的功能是进行文本的替换，替换的基本用法如下：

```
sed -e 's/pattern/replacement/options' filename
```

其中 *s* 表示“替换”，*pattern* 是要用来搜索替换文本的正则表达式，而 *replacement* 是要替换的结果，*filename* 是要操作的文件，*options* 是可选参数。如果必须同时执行多个替换操作，则每个操作之前都必须使用 `-e`。

```
sed -e 's/pattern1/replacement1/' -e 's/pattern2/replacement2/' filename
```

在默认情况下，sed 只会对每行文本进行一次替换，如果要替换所有的文本，则应该设定 *options* 参数为 *g*：

```
sed -e 's/pattern/replacement/g' filename
```

## 16.5.2 字符组及简记法

sed 的字符组遵循 POSIX 规范，所以 POSIX 的字符组和简记法都可以直接在 sed 中使用。需要指出的是，sed 中的点号 `.` 可以匹配包括换行符在内的任何字符，不过一般来说 sed 是按行处理的，所以这点区别并不重要。

从 sample.txt 中删除所有 a 或者 b

```
sed -e 's/[ab]//' sample.txt
```

从 sample.txt 中删除所有数字字符

```
sed -e 's/[[:digit:]]//g' sample.txt
sed -e 's/[0-9]//g' sample.txt
```

## 16.5.3 锚点和单词边界

sed 支持 `^` 和 `$`，其中 `^` 匹配行的起始位置，`$` 匹配行的结束位置。

从 sample.txt 中替换掉所有在行首的 cat

```
sed -e 's/^cat//g' sample.txt
```

从 sample.txt 中替换掉所有在行尾的 cat

```
sed -e 's/cat$//g' sample.txt
```

sed 认定的“单词字符”是 `\w`，也就是 `[[:alnum:]]`，在此基础上，共有 4 个单词边界：`\<`、`\>`、`\b`、`\B`。`\<` 匹配单词的起始位置，`\>` 匹配单词的结束位置，`\b` 匹配单词的起始或者结束位置，`\B` 则匹配 `\b` 无法匹配的位置。

从 sample.txt 中替换所有单词开头的 cat

```
sed -e 's/\<cat//g' sample.txt
```

从 sample.txt 中替换所有单词末尾的 cat

```
sed -e 's/cat\>//g' sample.txt
```

从 sample.txt 中替换所有的单词 cat

```
sed -e 's/\bcat\b//g' sample.txt
sed -e 's/\<cat\>//g' sample.txt
```

从 sample.txt 中替换所有单词内部的 cat（比如 truncate）

```
sed -e 's/\Bcat\B//g' sample.txt
```

### 16.5.4 量词

前面讲过，sed 属于 BRE 流派，所以虽然能使用的量词有 `*`、`+`、`?`、`{n, m}`、`{n}`、`{n, }`，但只有 `*` 可以直接使用，其他量词都必须转义：`+` 必须写成 `\+`，`?` 必须写成 `\?`，`{` 和 `}` 都需要转义。

从 sample.txt 中删除所有的 a、ab、abb、abbb...

```
sed -e 's/ab*//g' sample.txt
```

从 sample.txt 中删除所有的 ab、abb、abbb...

```
sed -e 's/ab\+//g' sample.txt
```

从 sample.txt 中删除所有的 a 或者 ab

```
sed -e 's/ab\?//g' sample.txt
```

从 sample.txt 中删除所有的 ab、abb、abbb

```
sed -e 's/ab\{1,3\}//g' sample.txt
```

### 16.5.5 多选结构和捕获分组

sed 中多选结构的 `[`、`|` 和 `]` 都需要转义，写成 `\( ... \)`。

从 sample.txt 中替换掉所有的 tom 或者 jerry

```
sed -e 's/\(tom|jerry\)//g' sample.txt
```

如果出现了捕获分组，则无论在 `regex` 中，还是 `replacement` 中，都可以用 `\1`、`\2` 的形式来引用。

从 sample.txt 中删掉重复的单词

```
sed -e 's/\(\b[[:alnum:]]\+\b\)[[:space:]]\+/\1/g' sample.txt
```

从 sample.txt 中把重复的单词替换为单个单词

```
sed -e 's/\(\b[[:alnum:]]\+\b\)[[:space:]]\+/\1/\1/g' sample.txt
```

sed 中有个特殊变量 `&`，它表示正则表达式匹配的文本。这样，我们可以给所有的数字字符串加粗：

```
sed -e 's/[[:digit:]]\+/&</b>/g' sample.txt
```

### 16.5.6 options

sed 的主要 *option* 介绍如下：

**g:** 在所有可能的地方进行替换

把 sample.txt 中所有的 tom 替换为 jerry

```
sed -e 's/tom/jerry/g'
```

**p:** 输出发生了替换的行。

它配合 -n 使用，可以只输出发生了替换的行。

把 sample.txt 中所有的 tom 替换为 jerry，并且只输出发生了替换的行

```
sed -n -e 's/tom/jerry/gp'
```

**l:** 查找-替换时不区分大小写（只有 GNU sed 支持此功能）

把 sample.txt 中所有的 tom（不区分大小写）替换为 jerry

```
sed -e 's/tom/jerry/Ig' sample.txt
```

**M:** 启用多行模式（只有 GNU sed 支持此功能）

在默认情况下，sed 按行读入文本，然后处理，但我们也可以指定 N 参数，让 sed 一次读入多行文本，此时如果指定了 M 模式，**^**和**\$**就能匹配多行文本中某一行的起始和结束位置。

一次读入两行，将行首的 tom 替换为 jerry（不要忘记 g，因为要处理“所有的行首”）

```
echo -e 'tom\ntom' | sed 'N; s/^tom/jerry/Mg
jerry
jerry
```

## 16.5.7 补充

如果觉得 BRE 的转义太麻烦，也可以指定使用 ERE 模式，办法就是在 sed 之后添加 -r 参数，如果之前使用了 -e 参数，可以把两个参数合并为 -re（写成 -er 会报错）。

**BRE**

```
sed -e 's/(\b[[:alnum:]]+\b)[[:space:]]+\1\b/\1/' sample.txt
```

**ERE**

```
sed -re 's/(\b[[:alnum:]]+\b)[[:space:]]+\1\b/\1/' sample.txt
```

sed 不支持环视，也不支持忽略优先量词。

关于 sed 的详细文档，可以输入 `info sed` 获得，也可以参考下面的文档：

<http://www.gnu.org/software/sed/manual/sed.html>

[http://www.linuxtopia.org/online\\_books/linux\\_tool\\_guides/the\\_sed\\_faq/index.html](http://www.linuxtopia.org/online_books/linux_tool_guides/the_sed_faq/index.html)

## 16.6   总结

本章介绍了 Linux/UNIX 下几款常用的文本编辑工具，它们的功能各不相同，但又有所重叠，为方便大家理解，表 16-11 列出了各款工具适合解决的问题类型。

表 16-11   几种工具的总结

工具	特性	适用问题
vi	全文编辑工具，通用文本编辑	全文编辑，微调整篇文本
grep	流编辑工具，不能修改文本，只能进行查找，结果中仍然保留原来行的形态	在大量文本中找出感兴趣的行，日志的筛选
awk	流编辑工具，不能修改文本，但结果不必保留原来行的形态，可以截取、重新组织各个字段，编程能力强大，可设定变量	筛选出感兴趣的字段，转换数据格式，生成统计报告
sed	流编辑工具，可以大幅度修改文本	批量修改大量文本（比如修改源代码中的链接），做临时修改配合 awk 和 grep 的使用

因为在实际开发中我们可能面对海量数据，用任何一款全文编辑器（包括 vi）打开编辑都是不现实的，只能用流编辑器来处理，而且必须混用 grep、awk、sed，让它们各司其职，配合完成任务。常见的流程是，先用 grep 筛选出真正感兴趣的行，再用 sed 修改数据、awk 整理格式并进行统计<sup>1</sup>。

下面给出一个真实的例子，作为本章的结尾。

某服务器记录了客户端软件的访问记录，其日志样本如下：

```
server1 - - [20/Dec/2010:06:33:58 +0800] "POST /get HTTP/1.1" 200 417 "-" "program-win/0.6.0.253 (b:i386; o:Microsoft Windows 7 Ultimate Edition (build 7600), 64-bit)" "219.70.179.147"
server1 - - [20/Dec/2010:06:33:59 +0800] "POST /get HTTP/1.1" 200 839 "-" "program-win/0.6.0.240 (b:i386; o:Microsoft Windows 7 Ultimate Edition Service Pack 1, v.721 (build 7601), 64-bit)" "122.70.32.72"
server1 - - [20/Dec/2010:06:33:59 +0800] "POST /get HTTP/1.1" 200 839 "-" "program-win/0.6.0.240 (b:i386; o:Microsoft Windows 7 Ultimate Edition Service Pack 1, v.721 (build 7601), 64-bit)" "122.70.32.72"
server1 - - [20/Dec/2010:07:13:27 +0800] "POST /get HTTP/1.1" 200 1078 "-" "program-linux/0.6.0.253 (X11; U; Linux x86_64; en-US; rv:1.9.2.13) " "125.114.169.249"
```

---

<sup>1</sup> awk 也可以执行 grep 的“筛选”功能，但 grep 可以智能混用正则表达式匹配和字符串的 Boyer-Moore 搜索算法，在进行简单处理时不需要动用正则表达式匹配，所以速度远远高于其他“纯”正则表达式工具。因此，进行文本筛选时 grep 是首选工具。

每一行记录一次访问，包含的信息有：所访问的服务器、日期、HTTP Header、时间、版本、客户平台、客户端 IP。

要完成的任务是统计 Windows 程序的版本分布状况，生成只包含版本号（不包含平台信息）和客户端 IP 地址的数据文件。上面的日志中有 3 行记录了 Windows 程序的访问，最后得到的数据文件如下：

```
0.6.0.253 219.70.179.147
0.6.0.240 122.70.32.72
0.6.0.240 122.70.32.72
```

一般来说，此类任务应当首先用 `grep` “筛选”出真正要处理的行（也就是包含 `program-win` 的行），然后用 `sed` 处理。常见的直接思路是用两个捕获分组匹配要提取的字段，将整行替换为这两个捕获分组的内容。整个程序如下：

```
grep 'program-win' log \|
sed -re 's/^.*'program-win\|/([^\ ]+).*"([^\"]+)"$/\1 \2/'
```

注意：这里用了 `-r` 选项，使用 ERE 写法，否则要增加不少元字符转义。

```
grep 'program-win' log \|
sed -re 's/^.*'program-win\|/([^\ ]+\|).*"([^\"]+\|)"$/\1 \2/'
```

思路是使用捕获分组抓取所需要的数据，将整行替换为这两个分组捕获的内容。结果没有错，但速度太慢，处理 200 万条数据需要 2 分钟左右，在生产环境中不可接受。

如果仔细思考就会发现，真正的问题只是“将感兴趣的内容提取出来，忽略其他内容”，捕获分组和替换并非必需。如果能直接把所需要的数据提取出来，就可以省去很多工夫。

但是，`awk` 的切分必须使用统一的分隔符，但要提取所需要的数据，有时需要以空白字符作为字段分隔，有时必须以 `/` 字符作为字段分隔，所以仅仅依靠 `awk` 是不行的。这时候，不妨动用 `sed`，将反斜线字符替换为空格，再切分。以 `grep` 提取出的第一行为例，原始数据是这样的：

```
server1 - - [20/Dec/2010:06:33:58 +0800] "POST /get HTTP/1.1" 200 417 "-" "program-win/
0.6.0.253 (b:i386; o:Microsoft Windows 7 Ultimate Edition (build 7600), 64-bit)"
"219.70.179.147"
```

先用双引号"作为分隔符，切分出完整的版本信息和客户端 IP 地址。

```
awk 'BEGIN { FS = "\"" } ; { print $6, $8 }'
```

```
program-win/0.6.0.253 (b:i386; o:Microsoft Windows 7 Ultimate Edition (build
7600), 64-bit) 219.70.179.147
```

这时候把反斜线替换为 `/` 空格，再切分，提取出第 2 个字段，以及最后一个字段（使用特殊变量 `$NF`）就可以完成：

```
sed -e 's/\\/ /g' | awk '{print $2, $NF}'  
0.6.0.253 219.70.179.147
```

最后把这些操作整合起来：

```
grep 'program-win' log \|  
awk 'BEGIN { FS = "\"" } ; { print $6, $8 }' \|  
sed -e 's/\\/ /g' | awk '{print $2, $NF}'
```

这样避免了使用捕获分组，而且大大减少了替换操作的工作量。经过测试，200 万条数据只用不到 10 秒钟就可以处理完成，速度提高了 10 倍以上。

从这个例子中，我们可以看到 **grep**、**awk**、**sed** 适合解决的问题：**grep** 适合过滤和筛选，**awk** 适合重新组织数据，**sed** 可以修改数据，也可以整理格式，为 **awk** 的处理提供便利。妥善结合运用这三种工具，可以大大提高解决问题的效率。

# 附录 A 常用语言中正则特性一览

字符及字符组						
特性	.NET	Java	JavaScript	PHP	Python	Ruby
\Q...\E	×	√	×	√	×	×
\d\w\s 字符组简记法	①	ASCII 匹 配规则	ASCII 匹 配规则	ASCII 匹配规则	②	ASCII 匹 配规则
POSIX 字符组	×	ASCII 字 符	×	ASCII 字符	×	③
<p>① 可以指定 <code>RegexOptions.ECMAScript</code> 恢复到 ASCII 匹配规则</p> <p>② Python 2 默认采用 ASCII 匹配规则，但指定 <code>Re.U</code> 模式可以变换为 Unicode 匹配规则；Python 3 默认采用 Unicode 匹配规则，但指定 <code>Re.A</code> 模式可以变换为 ASCII 匹配规则</p> <p>③ Ruby 1.9 的 POSIX 字符组支持 Unicode 字符，且不需要显式指定 Unicode 模式</p>						
括号相关						
特性	.NET	Java	JavaScript	PHP	Python	Ruby
表达式中 \1\9 引用分 组	√	√	√	√	√	√
替换中引用分组	<code>\$num</code>	<code>\$num</code>	<code>\$num</code>	<code>\$num</code>	<code>\g&lt;num&gt;</code>	<code>\num</code>
命名分组	④	×	×	⑤	⑥	⑦
<p>④ .NET 中用 <code>(?&lt;name&gt;regex)</code> 表示命名分组，用 <code>\k&lt;name&gt;</code> 在表达式中引用分组，用 <code>\${name}</code> 在替换中引用</p> <p>⑤ PHP 中用 <code>(?P&lt;name&gt;regex)</code> 表示命名分组，用 <code>(P=name)</code> 在表达式中引用分组，替换中不能引用命名分组</p> <p>⑥ Python 中用 <code>(?P&lt;name&gt;regex)</code> 表示命名分组，用 <code>(P=name)</code> 在表达式中引用分组，用 <code>\g&lt;name&gt;</code> 在替换中引用</p> <p>⑦ 只有 Ruby 1.9 支持命名分组，用 <code>(?&lt;name&gt;regex)</code> 表示命名分组，用 <code>\k&lt;name&gt;</code> 在表达式中引用分组，用 <code>\k&lt;name&gt;</code> 在替换中引用</p>						
断言						
特性	.NET	Java	JavaScript	PHP	Python	Ruby
\b (与 \w 规则相同)	√	⑧	√	√	√	⑨
^	√	√	√	√	√	⑩
\$	√	√	⑪	√	√	⑫





## 附录 B 常用的正则表达式

之前已经讲解了正则表达式的功能和用法，但在日常开发中，需要用正则表达式解决的各种问题经常是固定的、重复的，所以这里开列一些常用的正则表达式供查阅。

因为希望做的是尽量通用，所以只列出纯粹的正则表达式（为避免字符组简记法在不同语言中的规定不同，一般尽量不使用字符组简记法，如果确实出现了，除非特殊说明，否则一定是使用 **ASCII 匹配规则**，参见第 7 章），这个正则表达式在具体语言中要如何写（如果使用字符串，要如何转义），如何使用（调用哪个 API），都没有介绍，也没有刻意使用非捕获型括号 `(?:...)` 取代捕获型括号 `(...)`，如果你希望了解这些内容，请仔细阅读前面的章节。

另外，如果将正则表达式用于验证，一般应该在表达式首尾加上 `\A` 和 `\Z`；如果需要将其用于提取，应该在首尾分别加上对应的断言。这一点，除非特别要注意的地方会明确说明，一般不做说明。

### 整数

举例

1024

89

10

正则表达式

`[0-9]+`

如果用于数据提取，不需要添加断言。

### 十六进制数

举例

4e00

9F

A3D6

正则表达式

```
[0-9a-fA-F]+
```

如果用于数据提取，不需要添加断言。

## MD5 字符串

举例

```
C37B58C783456BF879C3DD5C519195A2
44912ED4AD09F0D2CA5066A16209F717
```

正则表达式

```
[0-9a-fA-F]{32}
```

常见的 MD5 字符串由 32 位十六进制数构成。

如果用于数据提取，还应该在首尾添加断言，保证之前和之后不能出现十六进制字符，也就是在表达式开头加上 `(?![0-9a-fA-F])`，在末尾加上 `(?![0-9a-fA-F])`。

## 浮点数

举例

```
+3.14
+0.07
.7
3.6
-3.14
-12
```

正则表达式

```
(+?([0-9]+|\.[0-9]+|\.[0-9]+\.[0-9]+)|-?([0-9]+|[0-9]+\.[0-9]+))
```

浮点数包括整数部分、小数点、小数部分，分别由 `\d+`、`\.`、`\d+` 匹配，单独来看它们都不一定要出现，但不可能三个部分都不出现，所以用多选结构 `([0-9]+|[0-9]+\.[0-9]+|\.[0-9]+)` 表示。另外之前可能出现符号，如果出现的是 `+`，则各种形式都没问题，如果出现的是 `-`，一般不会容许出现 `-.7` 之类的写法，所以还需要用一个多选结构。

如果用于数据提取，还应该在首尾添加断言。一般来说判断之前的字符不是点号，之后的字符一般不能是点号，但如果是句号，则它之后不能是其他数字字符，这样就可以避免错误匹配 `2010.12.20` 之类的数据。所以要在表达式开头加上 `(?!.)`，在末尾加上 `(?![0-9])`。

## 逗号分隔的整数

举例

123,456,789

23,456,789

3,456,789

正则表达式

```
(?<![0-9])[0-9]{1,3}(,[0-9]{3})*(?![0-9])
```

英文中的数字经常会加入逗号以便识别，规则是从右向左，每 3 位数字为一段，在之前加入一个逗号，比如 12,345,678。最左边一段数字的长度不能超过 3，所以表达式是 `[0-9]{1,3}`；其他每段长度都是 3，所以表达式是 `,[0-9]{3}`。

因为在这个表达式中必须对各段的长度做出限制，所以首尾应当分别使用断言 `(?<![0-9])` 和 `(?![0-9])`，确保之前和之后没有数字字符。如果你使用的是 Ruby 1.8，不支持逆序环视，也可以用 `\b` 替代 `(?<![0-9])`，作为权宜之计。

## 中文字符

在 Unicode 编码环境下（详见第 7 章）使用的正则表达式，匹配 Unicode 字符集中基本的中文字符（足够覆盖常见的绝大多数情况）。

**.NET Java JavaScript Python**

```
[\u4e00-\u9fff]
```

如果使用 Python 2 版本，必须在字符串之前添加 `u`，比如 `u"[\u4e00-\u9fff]"`。

**PHP**

```
[\u{4e00}-\u{9fff}]
```

必须指定 Unicode 模式，比如 `/[\u{4e00}-\u{9fff}]/u`。

**Ruby**

```
[\x{4e00}-\x{9fff}]
```

必须在 Ruby 1.9 以上版本中使用，而且必须显式指定 Unicode 模式，比如 `/[\x{4e00}-\x{9fff}]/u`。

在 GBK 编码环境下使用的正则表达式。

**Ruby 1.8 PHP Python**

```
[\xb0-\xfe][\x00-\xff]
```

这个正则表达式将中文字符作为两个字节加以匹配，所以不需要指定任何与 Unicode 有关的设置，无论是对正则表达式还是字符串，都是这样。而且，如果需要匹配多个中文字符，必须添加括号将两个字节作为一组，再以量词限定，比如 `([\xb0-\xfe][\x00-\xff])+`。

## 中文间多余的空白字符

举例

中文的 文本，往往会包含许多这样的空白字符  
但是，又不能简单将所有空白都删掉，否则 `you will get a big mess`

正则表达式

```
(?<=[\u4e00-\u9fff])\s+(?=[\u4e00-\u9fff])
(?<=[^\x00-\x7f])\s+(?=[^\x00-\x7f])
```

在填写中文内容时，因为转贴等操作，往往会多出许多空白字符，破坏格式，影响美观。但解决这类问题，又不能简单删掉所有空白字符，否则就破坏了其中的英文文本的格式，所以要删除的只是“中文字符之间的空白字符”（英文以空白字符作为词的分隔，而正常的中文文本之间很少出现空白字符）。如果要处理的文本是 Unicode 编码，则可以直接使用上面介绍的匹配中文字符的表达式 `[\u4e00-\u9fff]`，以逆序环视 `(?<=[\u4e00-\u9fff])` 和顺序环视 `(?=[\u4e00-\u9fff])` 要求匹配 `\s+` 的空白字符两侧都必须是中文字符，再将空白字符删除即可。如果要处理的文本不是 Unicode 编码，处理的思路仍然相同，只是假设所有的非 ASCII 字符都是中文字符，用排除型字符组 `[^\x00-\x7f]` 匹配，这个表达式有点风险，不过大多数情况还是够用的。

因为这个表达式就是用来提取（匹配-替换）的，所以不需要再添加断言。

## 价格

举例

```
¥189.99
$ 46
88.0
12700
```

正则表达式

```
[¥$]?[s*][0-9]+(\. [0-9]{2})?
```

价格比较好判断，之前可能有货币符号 ¥ 或者 \$，在 Unicode 编码环境中，¥ 是单个字符，可

以直接使用字符组 `[¥$]`，否则用 `(¥|\$)` 比较好，记得 `$` 必须转义，否则会作为锚点处理。在这个符号之后，具体的价格数值之间，可能还有若干空白字符，所以要使用 `\s*`。如果价格出现了小数部分，一般不会超过 2 位，所以用 `(\.[0-9]{1,2})`。

如果用于数据提取，还应该在尾部添加断言，判断之后的字符不是数字字符，也就是在表达式末尾加上 `(?![0-9])`。

## 手机号码

举例

```
1380123467
+861380123467
+86 1380123467
0086 1380123467
01380123467
1827654321
```

正则表达式

```
((0|(00|\+86)\s?)*(13[0-9]|15[0-356]|18[025-9])\d{8})
```

到本书写作时候为止，国内已经开放的手机号段有 130~139、150~153、155~156、180、182、185~189，用多选分支 `(13[0-9]|15[0-356]|18[025-9])` 可以很准确地匹配号段；之后的 8 位一般没有限制，只要是数字即可，用 `\d{8}` 匹配。另一方面，手机号之前可能有 0 或者 +86 或者 0086，而且如果出现，在它和后 11 位之间，可能还有一个空白字符，所以用 `(0|(00|\+86\s?))` 匹配，因为这部分是可能出现的，所以还需要加上量词，变为 `(0|(00|\+86\s?))?`。

如果用于数据提取，还应该在首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## 固定电话号码

举例

```
1234567
021-12345678
0731-12345678
(010)12345678
(010)12345678-836016
```

正则表达式

```
(\\(?0[1-9]{2,3}\\)?-?)?[1-9][0-9]\\{6,7\\}(-[0-9]{1,6})?
```

国内的电话号码分为两段，区号一般为 3~4 位，且以 0 开头，之后不为 0，所以是 `0[1-9]{2,3}`；市内号码一般为 7~8 位，且不能以 0 开头，所以是 `[1-9][0-9]\\{6,7\\}`；如果有分机，分机号一般不超过 6 位，且以 - 与总机号码相连，所以是 `-[0-9]{1,6}`。因为区号可能以括号标注（`(010)12345678`），也可能以 - 与市内号码相连（`010-12345678`），或者直接相连（`01012345678`），所以区号部分的表达式是 `(\\(?0[1-9]{2,3}\\)?-?)?`；分机号码也不一定出现，所以需要 `?` 限定，记为 `(-[0-9]{1,6})?`。

如果用于数据提取，还应该为首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## 邮政编码

举例

100859

412000

正则表达式

```
[1-9][0-9]{5}
```

邮政编码一般是 6 位数字，首位不能为 0，所以是 `[1-9][0-9]{5}`。

这个表达式比较简单，许多数据中都可能出现这样的字符串，所以提取时一定要不要忘记在两端添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## 身份证号码

举例

310100198001012071

10010019800101207x

310100800101207

正则表达式

```
[1-9][0-9]{14}([0-9]{2}[0-9xX])?
```

我国身份证分为两代，第一代身份证号码长为 15 位，第二代身份证号码长为 18 位，都不能以 0 开头，如果是 18 位身份证号，则最后一个字符可以是 x 或者 X（在特点情况下，按照公安系

统和网银系统的要求，只能使用大写 X)，其他都只能为数字。所以可以这么看，第一位是[1-9]，之后是 14 位[0-9]，如果是 18 位号码，则最后还会出现[0-9]{2}[0-9x]。

如果用于数据提取，还应该为首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上(?:<![0-9])，在末尾加上(?:![0-9])。

## QQ 号

举例

10000

95381376

正则表达式

[1-9][0-9]{4,9}

目前 QQ 号最长为 10 位，最短为 5 位，且首位不能为 0。

如果用于数据提取，还应该为首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上(?:<![0-9])，在末尾加上(?:![0-9])。

## ISBN 号

举例

ISBN 978-7-534-74286-6

ISBN-13 9787534742866

ISBN-13: 978-7-534-74286-6

9787534742866

0-596-52068-9

ISBN-10 0-596-52068-9

正则表达式

((ISBN(-13)??:?\s)?97[89][-\s]?[0-9][-\s]?[0-9]{3}[-\s]?[0-9]{5}[-\s]?[0-9]|(ISBN(-10)??:?\s)?[0-9][-\s]?[0-9]{3}[-\s]?[0-9]{5}[-\s]?[0-9x])

ISBN 有 10 位 (ISBN-10) 和 13 位 (ISBN-13) 两种长度，两者的主要区别是：ISBN-10 的开头可能是 ISBN-10，而 ISBN-13 的开头只可能是 ISBN-13；最后一位校验码，ISBN-10 中的取值可能是[0-9x]，而 ISBN-13 中只可能是[0-9]；ISBN-13 的开头 3 位只能是 978 或者 979。无论 ISBN10 还是 ISBN-13，最后的 10 位都可以按照 1、3、5、1 分为四段，各段之间可能以-或者空白字符作为分隔符，比如 7-534-74286-6 或者 7 534 74286 6，或者也可以没有任何分隔符，



比如 7534742866。

对应开头部分的正则表达式比较容易，ISBN-10 是 `(ISBN(-10)??:?\s)?`，ISBN-13 则是 `(ISBN(-13)??:?\s)?`。之后，如果是 ISBN-13，必然会出现 `97[89]`，再将最后 10 位按照 1、3、5、1 四段分开，中间加上可能出现的分隔符 `[-\s]?`。如果是 ISBN-10，则最后一位是 `[0-9x]`，ISBN-13 则为 `[0-9]`，为避免两种情况混淆，将 ISBN-10 和 ISBN-13 两种情况分别作为一个多选分支，列在同一个多选结构内。

如果用于数据提取，还应该在首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## IPv4 地址

举例

```
192.168.1.1
202.198.0.68
```

正则表达式

```
(0{0,2}[0-9]|0?[0-9]{2}|1[0-9][0-9]|2[0-4][0-9]|25[0-5]\.){3}(0{0,2}[0-9]|0?[0-9]{2}|1[0-9][0-9]|2[0-4][0-9]|25[0-5])
```

这个例子在第 3 章讲解过，IPv4 的每一段都在 0~255 之间，为了考虑一位数、两位数、三位数的情况，必须用多选结构来解决。注意这个表达式中不能使用反向引用，而应当使用量词。

如果用于数据提取，还应当在首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## IPv6 地址

举例

```
fe80:0:0:09a:fe:0:0:4ca2
da8:207:3:21a:64ff:fe6d:d1eb
```

正则表达式

```
([0-9a-fA-F]{1,4}:){7}[0-9a-fA-F]{1,4}
```

IPv6 地址一般分为 8 段，每段都用十六进制数值表示，长度在 1~4 之间，所以是 `[0-9a-fA-F]{1,4}`。

如果用于数据提取，还应当在首尾添加断言。一般来说，判断之前和之后的字符不是数字字

符即可，也就是在表达式开头加上 `(?<![0-9a-fA-F])`，在末尾加上 `(?![0-9a-fA-F])`。

## 时间字符串

24 小时制，如 20:46

正则表达式

```
(0?[0-9]|1[0-9]|2[0-3]):(0?[0-9]|1[234][0-9]|5[0-9])
```

12 小时制，如 12:38

正则表达式

```
(0?[0-9]|1[012]):(0?[0-9]|1[234][0-9]|5[0-9])
```

时间字符串的匹配与之前 IPv4 地址的匹配类似，需要注意的是数据的取值范围：如果是 24 小时制，小时数只能在 0~23 之间，所以正则表达式是 `(0?[0-9]|1[0-9]|2[0-3])`；如果是 12 小时制，小时数则只能在 0~12 之间，正则表达式是 `(0?[0-9]|1[012])`；分钟数只能在 0~59 之间，所以正则表达式是 `(0?[0-9]|1[234][0-9]|5[0-9])`。如果还要包括秒数，则还应该附加 `:(0?[0-9]|1[234][0-9]|5[0-9])`。

如果用于数据提取，还应当在首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## 日期字符串

以 mm/dd/yyyy 为例，如 12/20/2010

正则表达式

```
(0?[1-9]|1[012])/(0?[1-9]|1[12][0-9]|3[01])/[0-9]{4}
```

时间字符串的匹配与之前 IPv4 地址的匹配类似，需要注意的是数据的取值范围：月只能在 0~12 之间，所以正则表达式是 `(0?[1-9]|1[012])`；日只能在 0~31 之间（暂不考虑大、小月），所以正则表达式是 `(0?[1-9]|1[12][0-9]|3[01])`；年一般无限制，所以正则表达式是 `[0-9]{4}`。

如果用于数据提取，还应当在首尾添加断言。一般来说，判断之前和之后的字符不是数字字符即可，也就是在表达式开头加上 `(?<![0-9])`，在末尾加上 `(?![0-9])`。

## Windows 路径

举例

```
C:\windows
e:\Program Files\
d:\starcraft2\startcraft.exe
```

正则表达式

```
[a-zA-Z]:\\([^\\"/:*<>|"?\\r\\n]+\\)*[^\\"/:*<>|"?\\r\\n]+\\/?
```

Windows 下的路径，开头是盘符，必须是一个英文字母，然后是:，注意\在正则表达式中必须转义为\\，匹配这部分的正则表达式是[\[a-zA-Z\]:\\](#)；文件名（包括目录名）中不能包含 \ / : \* < > | " ? 字符，还不能是换行符，也不能是空字符串，所以表达式是[\[^\\"/:\\*<>|"?\\r\\n\]+](#)，目录之间以\分隔，路径内部的目录名之后必须有\，所以这些部分可以用[\(\[^\\"/:\\*<>|"?\\r\\n\]+\\\)\\*](#)匹配；路径末尾的文件名如果是目录，可能会以\结尾，所以最后的部分是[\[^\\"/:\\*<>|"?\\r\\n\]+\\/?](#)。

如果用于数据提取，还应当为首尾添加断言，在表达式开头加上[\(?:!\[a-z\]\)](#)，在末尾加上[\(?:!\[^\S\\\"/:\\*<>|"?\]\)](#)即可。

## UNIX 路径

举例

```
/usr/local
/usr/local/
/usr/bin/python
bin/python
```

正则表达式

```
/(?:[^\\"/:*%<>|"?\\r\\n]+)/([^\\"/:*%<>|"?\\r\\n]+)/?
```

按照 POSIX 规范，文件名不能包含 \ / : \* % < > | " ? 字符，不能是换行符，也不能是空字符串，所以匹配文件名的表达式是[\[^\\"/:\\*%<>|"?\\r\\n\]+](#)，在路径内部的所有目录名之后必须有/，所以这些部分可以用[\[^\\"/:\\*%<>|"?\\r\\n\]+/](#)匹配；在路径的开头可能有/，所以要加上[/?](#)；路径末尾的文件名如果是目录，可能以/结尾，所以最后的部分是[\[^\\"/:\\*%<>|"?\\r\\n\]+/?](#)。

如果用于数据提取，还应当为首尾添加断言，一般来说，在表达式开头加上[\(?:!\[^\S\\\"/:\\*%<>|"?\\r\\n\]\)](#)，在末尾加上[\(?:!\[^\S\\\"/:\\*%<>|"?\\r\\n\]\)](#)。

## 主机名

举例

```
localhost
163.com
news.163.com
sina.com.cn
```

正则表达式

```
(?=[-a-zA-Z0-9.]{0,255}(?![-a-zA-Z0-9.]))((?!-)[-a-zA-Z0-9]{1,63}\.)*((?!-)[-a-zA-Z0-9]{1,63})
```

主机名分为很多个字段 (label)，每个字段可以包含字母、数字、横线，其长度不能超过 63，所以是 `[-a-zA-Z0-9]{1,63}`；且字段不能以横线开头，这一点用 `(?!-)` 来保证。总长度不能超过 255，用 `(?=[-a-zA-Z0-9.]{0,255}(?![-a-zA-Z0-9.]))` 来保证。

整个主机名可以看做“字段+点号”的重复，所以是 `((?!-)[-a-zA-Z0-9]{1,63}\.)*`，最后的字段必须要出现，所以最后再添加一个匹配字段的子表达式。

如果用于数据提取，则必须在表达式开头添加断言 `(?<![a-zA-Z0-9.])`。

## 电子邮件地址

举例

```
someone@localhost
tom@163.com
jerry.lee@sina.com.cn
```

正则表达式

```
(?!.)(?![\w.]*?\.\.)([\w.]{1,64})@(?=[-a-zA-Z0-9.]{0,255}(?![-a-zA-Z0-9.]))((?!-)[-a-zA-Z0-9]{1,63}\.)*((?!-)[-a-zA-Z0-9]{1,63})
```

电子邮件地址由用户名和主机名两部分构成，以@连接，主机名的匹配之前已经讲过，这里只讲解用户名的匹配。

一般来说，用户名中能出现的字符是字母、数字、下画线、点号，其长度不超过 64 个字符，所以是 `[\w.]{1,64}`；用户名不能以点号开头，所以要添加环视 `(?!.)`，同时用户名中不能包含连续点号，所以还需要添加环视 `(?![\w.]*?\.\.)`。

如果用于数据提取，还必须在之前添加断言 `(?<![\w.])`。

## URL

举例

```
http://www.yahoo.com
http://www.yahoo.com/
http://yahoo.com:80
http://news.163.com
http://news.163.com/11/0625/03/77C5GLMJ0001124J.html
http://baike.baidu.com/view/1203256.htm
http://www.baidu.com/s?wd=%D5%FD%D4%F2
```

正则表达式

```
(https?|ftp)://[^\?:]�+(:[0-9]{1,5})?(/?(/[^\?]*)(\?[\^s"'']�)?)
```

相对来说，匹配 URL 的表达式比较复杂，因为关于主机名（hostname）、路径名（path）、参数（parameter）等部分都有复杂的规范（所以不少语言中都提供了专门的 URLParse 函数），这里只给出一个相对简单的正则表达式，在日常使用中可以满足需求。

最开始是协议名，可能是 http、https、ftp，这部分对应的表达式是 `(https?|ftp)`，然后是 `://`；主机名使用了相对简单的 `[^\?:]�` 匹配（完整的匹配可以参考匹配之前匹配主机名的表达式），然后用 `(:[0-9]{1,5})?` 匹配可能出现的端口部分（因为端口号最大为 65536，这里简要规定为 5 位数字）。如果这个 URL 只包含主机名（以及可能的端口号），可能到此为止，也可能之后还有 `/`（比如 `http://www.yahoo.com/`）；如果之后还有内容，一般是以 `/` 分隔开的路径名，所以用 `(/[^\?]*)*` 匹配，最后可能还包括参数（比如 `?wd=%D5%FD%D4%F2`），用 `(\?[\^s"'']�)?` 匹配（这里假定 URL 中不包含空字符和引号字符，而且 URL 之后必然有空字符，或者位于字符串的末尾）

URL 的规范比较多，这个表达式能进行一般的验证，也能在相对规范（通常的文档）的数据中提取出 URL，如果存在很多噪音数据（比如似是而非的 URL），这个表达式可能无能为力，此时请参考其他资料。

## HTML Entity

举例

```
&#x53d1;
&#21457;
```

正则表达式

```
&#([0-9]+|x[0-9a-fA-F]�);
```

HTML Entity 是 HTML 代码中常见的数据，格式为 `&#xhex;` 或者 `&#dec;`，其中的 *hex* 或者 *dec* 为所示字符的码值，*hex* 为十六进制，*dec* 为十进制，所以 `&#x53d1;` 和 `&#21457;` 其实都表示 Unicode 编码的“发”。表达式开头的 `&#` 和结尾的 `;` 都不难理解，其中的多选结构中，`[0-9]+` 和 `x[0-9a-fA-F]+` 分别用来匹配十进制和十六进制数值字符串。

如果用于数据提取，不需要添加断言。

## HTML tag

举例

```
<a href="http://www.yahoo.com">

</h1>
```

正则表达式

```
(?i)</?[a-z][-a-z0-9_:.]*(?=[\s>))('^[']*'|"^["]*"|'>)*>
```

按照规范，HTML tag 必须有 tag name，第一个字符（close tag 中是 / 之后的字符，为了应对 close tag，在开头加上 `/?`）必须是字母 `[a-z]`，其中可以包含的字符还有数字 `[0-9]`、横线 -、下画线 \_、冒号 :、句号 .，所以用 `[a-z][-a-z0-9_:.]*` 匹配 tag name，为保证它完整匹配 tag name，用肯定顺序环视 `(?=[\s>])` 保证它之后是空白字符或者 `>`。在结尾的 `>` 之前可能出现的文本，只有在单引号字符串 `'^[']*'` 或双引号字符串 `"^["]*"` 内部可能出现 `>`，否则不能出现 `>`，这部分内容用 `'>'` 匹配，最后是结尾的 `>`。同时，这个表达式指定了不区分大小写模式 `(?i)`，保证 tag name 可以为大写。

如果用于数据提取，不用添加断言。

## Title tag

举例

```
<title>标题</title>
<TITEL>网站标题</TITEL>
<Title>标题</Title>
```

正则表达式

```
(?i)<title>.*?</title>
```

这个表达式比较简单，不多解释，只是注意使用不区分大小写模式，应对各种可能的情况。

## 成对的 tag

举例

```
<h1>yahoo</h1>
<span>yahoo</span>
<div></div>
```

正则表达式

```
(?is)<([a-z][-a-z0-9_:.]*)?(?=[\s>])(('[^']*'|"([^"]*"|'[^']*>))*).*?</\1>
```

这个正则表达式与上一个很类似，只是给匹配 tag name 的表达式加上一个捕获型括号，并在表达式结尾以 `</\1>` 引用对应的 close tag。两个 tag 之间的内容用 `. * ?` 匹配，注意必须使用忽略优先量词 `* ?`，并且必须使用单行模式（为了照顾 tag name 的大写情况，同时指定了不区分大小写模式）`(?is)`。因为 JavaScript 不支持单行模式，所以在 JavaScript 中必须将 `. * ?` 改为 `[\s\S]*?`，将 `(?is)` 改为 `(?i)`。

如果用于数据提取，不用添加断言。

## img tag

举例

```

<img src='pic.jpg' />
<img src = pic.jpg />

```

正则表达式

```
(?i)<img\s('[^']*'|"([^"]*"|'[^']*>))*src\s*=\s*['"]?[^'"\\s]+['"]?(['^']*'|"([^"]*"|'[^']*>))*/*>
```

img tag 以 img 开头，而且之后至少有一个空白字符，用 `\s` 匹配，然后用之前讲过的 `(('[^']*'|"([^"]*"|'[^']*>))*` 匹配 tag 内 src 之前可能出现的内容，用 `src\s*=\s*['"]?[^'"\\s]+['"]?` 匹配图片的地址，最后再用 `(('[^']*'|"([^"]*"|'[^']*>))*` 匹配其他可能的内容，最后是 `/>`。注意这个表达式同样要使用 `(?i)` 指定不区分大小写模式。

如果用于数据提取，可以把 `src\s*=\s*['"]?[^'"\\s]+['"]?` 改为 `src\s*=\s*['"]?([^'"\\s]+)['"]?`，即加上捕获型括号，匹配成功之后提取编号为 2 的分组即可。

## table tag

举例

```
<table><tr><td>1</td></tr></table>
<table border="1"><tr><td>1</td></tr></table>
```

正则表达式

```
(?is)<table\s(?:=\s[>])?(['^"]*"|'["^"]*"|'["^"]">])*.*?</table>
```

table tag 以 tab 开头，之后可能是空白字符也可能是>，所以用 `(?=[\s>])`，然后用之前讲过的 `(['^"]*"|'["^"]*"|'["^"]">)*` 匹配 tag 内可能出现的内容，再往后是>。close tag 用 `</table>` 匹配，之间的内容用 `.*` 匹配，注意必须使用忽略优先量词 `.*?`，并且必须使用单行模式（为了照顾 tag name 的大写情况，同时指定了不区分大小写模式）`(?is)`。因为 JavaScript 不支持单行模式，所以在 JavaScript 中必须将 `.*` 改为 `[\s\S]*?`。

如果用于数据提取，一般还需要给 `.*` 加上捕获型括号变为 `(..*?)`，匹配成功之后用分组 2 提取出 table 的内容，再按照 tr tag 拆分成行，每一行按照 td 拆分成列，得到单元格的内容，匹配 tr tag 和 td tag 的正则表达式与 table tag 的非常相似，此处不重复。

## href tag

举例

```
<a href="http://www.yahoo.com">yahoo</a>
<a href='http://www.yahoo.com'>yahoo</a>
<a href = http://www.yahoo.com>yahoo</a>
<a id="1" href="http://www.yahoo.com" class="do-edit">yahoo</a>
<a href="http://www.yahoo.com"></a>
```

正则表达式

```
(?i)<a\s(['^"]*"|'["^"]*"|'["^"]">)*href\s*=\s*['"]?[^'\s]+['"]?(['^"]*"|'["^"]*"|'["^"]">)*>.*?</a>
```

这个表达式类似匹配 img tag 的表达式，只是将 src 属性改为 href 属性，另外加上了匹配 close tag 的 `</a>`。在 open tag 和 close tag 之间的内容用 `.+?` 匹配，注意不能用 `[^<]+?`，因为其中还有其他 tag，比如 `<a ...><img ...></a>`，另外要使用忽略优先量词 `.*?`，防止错误匹配。

如果用于数据提取，不用添加断言，但是应该给 `['^"]*\s+` 和 `.+?` 分别加上捕获型括号，匹配成功之后用分组 1 提取链接的地址，用分组 3 提取链接的文本。

## JavaScript 代码

举例

```
<script type="text/javascript">alert("1")</script>
```



```
<script type="applicate/javascript">alert("1")</script>
<script language="text/javascript">alert("1")</script>
<script>alert("1")</script>
```

正则表达式

```
(?is)<script[\s].+?</script>
```

JavaScript 代码在<script 之后可能并不是空格, 而是空白字符, 再之后, 可能是 type="text/javascript", 也可能是 type="application/javascript", 还可能用 language 取代 type (实际上 language 是以前的写法, 现在大都用 type), 甚至可能没有属性, 直接是<script><sup>1</sup>。所以这个表达式应该改造, 不妨将条件放宽, 在 script 之后, 可能出现空白字符, 也可能直接是>, 可以用一个字符组 [\s>] 来匹配, 之后的内容统一用 .+? 匹配, 忽略优先量词保证了它会匹配到 </script> 结束。最终得到的表达式就是 <script[\s>].+?</script>。注意这个表达式要用 (?is) 指定不区分大小写模式和单行模式, 前者保证可以兼容<SCRIPT>的情况, 后者保证 .+? 匹配的文本可以跨越多行。因为 JavaScript 语言中的正则表达式不支持单行模式, 所以如果用在 JavaScript 程序中, 应当将 .+? 改为 [\s\S]+?, 将 (?is) 改为 (?i)。

与上面几个匹配 tag 的表达式相比, 似乎不够严谨, 但实际使用中并不会有很大的问题, 因为提取各种 tag 时, 更感兴趣的是 open tag 和 close tag 之内的内容而不是 tag 本身, 而提取 JavaScript 代码时, 一般提取整段即可。

如果用于数据提取, 不用添加断言。

## XML 格式注释

举例

```
<!--some -->
```

正则表达式

```
(?s)<!--.*?-->
```

XML 文档中的注释以<!--开头, 以-->结束。因为 JavaScript 不支持单行模式, 所以在 JavaScript 中必须将 .\*? 改为 [\s\S]\*?, 将 (?is) 改为 (?i)。

如果用于数据提取, 不用添加断言。

---

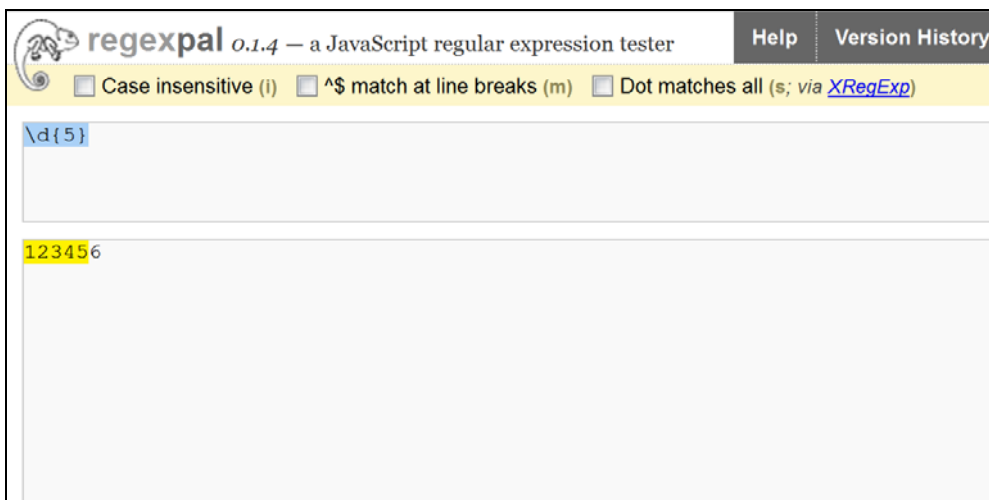
<sup>1</sup> 严格说起来, 如果只出现<script>, 无法保证这里出现的就是 JavaScript 代码, 也可能是 VBScript 代码。但考虑到真实世界中的情况, 可以认为<script 标识的“就是”JavaScript 代码, 所以这里不作区分。

## 附录 C 常用的正则表达式工具及资源

### 在线类

#### RegexPal

<http://www.regexpal.com>



这个网站可以测试 JavaScript 的正则表达式，在上面的文本框中输入正则表达式（不需要带上两端的分隔符），在下面的文本框中输入要处理的文本，运行结果会以高亮文字实时标注出来；各个匹配模式列在最上端，可以勾选。它主要进行正则表达式的验证和匹配，不能进行替换和切分。

#### Larsolavtorvik

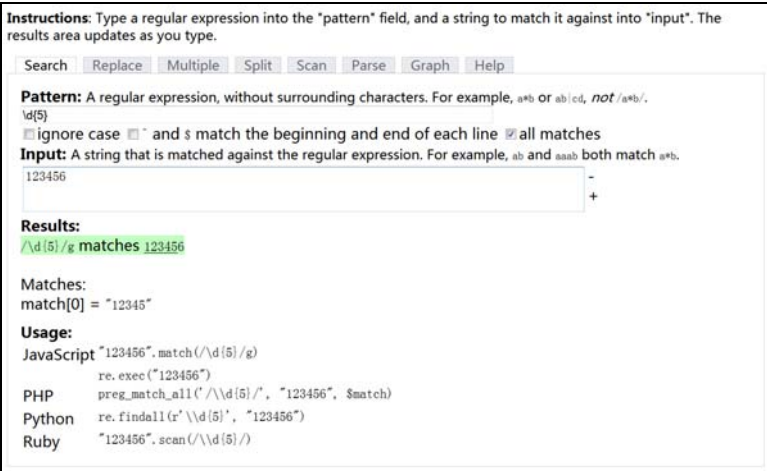
<http://regex.larsolavtorvik.com/>



这个网站用来调试 PHP 和 JavaScript 的正则表达式，它的功能比 [regexpal.com](http://regexpal.com) 强大很多。以 PHP 为例，先在左上角选择 PHP PCRE（常用的 preg 系列函数也就是默认的 PHP PCRE 选项），之后页面上会出现若干文本框，在 Pattern 中输入正则表达式（不需要带上两端的分隔符），在 Replacement 中可以输入 replacement 字符串，在 Subject 中输入要处理的文本，在 Code 栏中会自动生成 PHP 代码，在 Result 栏中会显示程序的运行结果。在右侧提供了各种函数、匹配模式、flags 的选择，可以勾选。

## ReWork

<http://osteele.com/tools/rework/>



这个网站同时提供了 JavaScript、Java、Python、Ruby 的正则表达式调试功能，在其中可以进行查找、替换、切分等多种操作，并且在 Mutiple 中可以同时显示对多个字符串的处理结果，

方便对比（但是它不支持环视）。在 **Pattern** 栏输入正则表达式，在 **Input** 栏输入要处理的字符串，可以在 **Pattern** 栏下勾选匹配模式对应的选项。之后会在 **Result** 栏实时用高亮标注匹配结果。

## NRegex

<http://nregex.com/nregex/default.aspx>

Nregex

How to use Nregex

Show Help

Show Useful

Feedback

> Nregex Bookmarklet

Regular Expression:

☐ Manually evaluate regex (large documents, latency issues)

☐ Ignore Case

☐ Single Line

☐ Multi Line

☐ Explicit Capture

Replacement String:

Matches & Replacements

C#

Load Target From URL:

Load

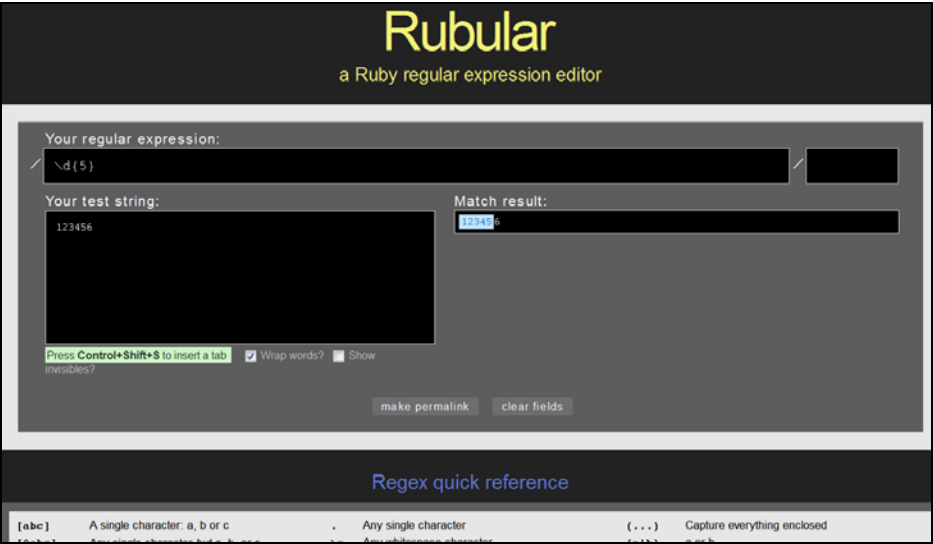
Load Target From File:

If I just had \$5.00 then "she" wouldn't be so @#\$.! mad.

这个网站用来调试.NET (C#) 中的正则表达式。在 **Regular Expression** 栏输入正则表达式（如果要设定匹配模式，可以勾选下方对应的选项），在最下方的文本框中输入要处理的字符串（也可以选择从文件或者 URL 导入）；如果要进行替换，可以在 **Replacement String** 中输入 replacement 字符串。之后，就会在 **Matches & Replacements** 栏显示程序运行的结果，它右侧的 **Matched Groups** 中详细显示了匹配的细节。

## Rubular

<http://www.rubular.com>

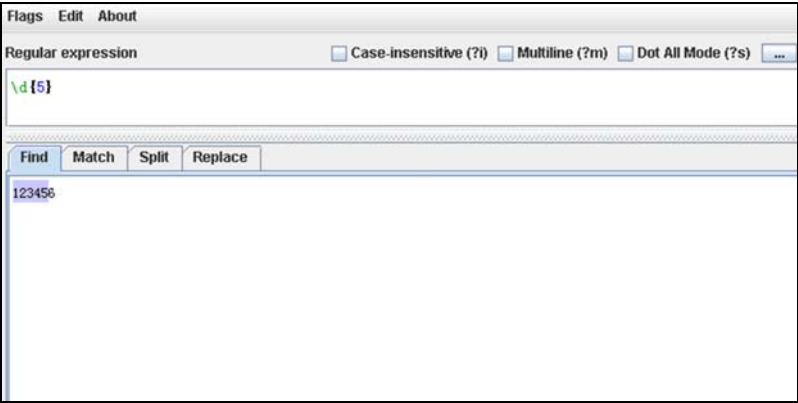


这个网站用来调试 Ruby 1.8 中的正则表达式。它非常形象，最上面的两个文本框，即两个分隔符/之间的文本框用来输入正则表达式，第二个分隔符/右侧的文本框用来指定匹配模式；将要处理的文本输入 Your test string 文本框，右侧的 Match result 中会用高亮标注出匹配结果。

因为这个网站只支持 Ruby 1.8，所以许多 Ruby 1.9 才提供的功能（比如逆序环视）都不能使用。

## MyRegexp

<http://myregexp.com/>



这个网站用来调试 Java 中的正则表达式（使用这个网站的调试功能，需要使用 Applet，要求本机有 JRE）。在 Regular expression 栏输入正则表达式，如果需要设定匹配模式，再选择下面的

Find、Match、Split、Replace 四个选项卡，可以进行对应的操作。如果要设定匹配模式，可以勾选右侧的选项，也可以从 **Flags** 下拉菜单中选择；**Edit** 下拉菜单中提供了几个很有用的功能：将当前表达式转换为 **Java 源代码**、**XML 源代码**、**JavaScript 源代码**，也可以从 **Java 源代码**、**XML 源代码**、**JavaScript 源代码** 获得正则表达式，这样就省去了繁杂的转义。

## RegexLib

<http://www.regexlib.com/>



这个网站堪称正则表达式的知识大全。它内含了许多关于正则表达式的讨论和总结，首页的 **Search** 功能可以找到与各种任务（比如 E-mail、URL、password）相关的正则表达式，**Browse Expressions** 则可以浏览网站内总结的表达式，并且提供 **Email/Uri/Numbers/Strings/Dates and Times** 等诸多分类；另外值得一提的是 **Regex Tester**，其中提供了针对 .NET 和 Silverlight 的正则表达式测试，在其中可以选择各种正则引擎（.NET、Client-side、Silverlight），在 **Source** 栏输入要处理的文本（也可以从外部文件或者 URL 导入），在 **Regular Expression** 栏输入正则表达式，点击 **Submit** 按钮，则会出现匹配结果。

## Regular-Expressions

<http://www.regular-expressions.info/>



这个网站同样囊括了与正则表达式相关的各种知识，它比 [regexlib.com](http://regexlib.com) 更加全面，其中不但有知识点的讲解，还有各语言的实现细节参考列表，并给出了许多例子，非常适合参考和引用。

## 软件类

### Expresso

<http://www.ultrapico.com/Expresso.htm>

这是 .NET 下的一个正则表达式测试工具，它能显示匹配的细节，方便确定问题所在，推荐使用。这个软件需要注册，但注册是免费的。

### RegexBuddy

<http://www.regexbuddy.com>

这是 Windows 下关于正则表达式的很完善的工具，它可以很方便地调试正则表达式，显示各部分匹配的细节。依靠这些功能，很容易理解复杂的表达式。它还可以自动为各种编程语言生成对应的代码片段，非常实用。这个软件需要注册使用，收费 39.95 美元。

### PowerGREP

<http://www.powergrep.com>

PowerGREP 堪称 Windows 下最全面、最强大的正则表达式工具，不但可以调试各种正则表达式，还可以用正则表达式完成各种文本编辑、文件改名等任务。这个软件需要注册使用，收费 159 美元。

## JFLAP

<http://www.cs.duke.edu/cseds/jflap/>

JFLAP 是杜克大学计算机系提供的，用来学习形式语言及自动机理论的工具，可以很方便地创建 DFA、NFA、正则语法、正则表达式，也可以进行正规语言、正则表达式、自动机之间的转换，最重要的是，它可以很形象地以图形说明原理和过程。JFLAP 以 jar 文件的方式提供，运行它要求本地有 JDK。

## RegexUtil

<http://sourceforge.net/projects/regex-util/>

这个软件是之前提到的 myregexp.com 提供的，它是 Eclipse 的一个插件，用来调试 Java 的正则表达式。安装之后，在 Eclipse 中选择 Window→Show View→Other→Regex Util 即可运行。